



Machine learning to predict impulse control disorders in Parkinson's disease

Johann Faouzi

► To cite this version:

Johann Faouzi. Machine learning to predict impulse control disorders in Parkinson's disease. Artificial Intelligence [cs.AI]. Sorbonne Université, 2020. English. NNT : 2020SORUS048 . tel-03090079v2

HAL Id: tel-03090079

<https://hal.science/tel-03090079v2>

Submitted on 26 Mar 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



SORBONNE UNIVERSITÉ

École Doctorale d'Informatique, de Télécommunication et d'Électronique (EDITE)

Laboratoire ICM – Équipe de recherche ARAMIS

Machine Learning to Predict Impulse Control Disorders in Parkinson's Disease

Johann FAOUZI

Thèse de doctorat d'informatique

Dirigée par Olivier COLLIOT et Jean-Christophe CORVOL
et co-encadrée par Baptiste COUVY-DUCHESNE

Présentée et soutenue publiquement le 4 décembre 2020

Composition du jury :

<i>Rapporteurs :</i>	Christophe AMBROISE (Professeur) Franck DURIF (PU-PH)
<i>Examineurs :</i>	Marta AVALOS (Maître de conférences) Gérard BIAU (Professeur) Alexandre GRAMFORT (Directeur de recherche) Mélanie PRAGUE (Chargée de recherche)
<i>Directeurs de thèse :</i>	Olivier COLLIOT (Directeur de recherche) Jean-Christophe CORVOL (PU-PH)

Abstract

MACHINE LEARNING TO PREDICT IMPULSE CONTROL DISORDERS IN PARKINSON’S DISEASE

by JOHANN FAOUZI

Impulse control disorders are a class of psychiatric disorders characterized by impulsivity. These disorders are common during the course of Parkinson’s disease, decrease the quality of life of subjects, and increase caregiver burden. Being able to predict which individuals are at higher risk of developing these disorders and when is of high importance.

The objective of this thesis is to study impulse control disorders in Parkinson’s disease from the statistical and machine learning points of view, and can be divided into two parts. The first part consists in investigating the predictive performance of the altogether factors associated with these disorders in the literature. The second part consists in studying the association and the usefulness of other factors, in particular genetic data, to improve the predictive performance.

In the first chapter, we present Parkinson’s disease and impulse control disorders, review the literature on impulse control disorders in Parkinson’s disease, introduce the main concepts of machine learning, and describe the databases from which we obtained data and the software used to analyze these data. In the second chapter, we investigate the predictive performance of several machine learning algorithms using features that have been associated with impulse control disorders in Parkinson’s disease. In the third chapter, we investigate the association between impulse control disorders in Parkinson’s disease and genetic risk scores for a broad range of phenotypes. In the last chapter, we investigate different approaches to integrate static data in recurrent neural networks and evaluate their predictive performance in the use case of predicting impulse control disorders in Parkinson’s disease, with genetic data used as static data.

Across these works, we highlight the importance of using machine learning algorithms, cross-validation and replication cohorts to unbiasedly estimate the predictive power of known and putative risk factors of impulse control disorders in Parkinson’s disease.

Résumé

APPRENTISSAGE AUTOMATIQUE POUR LA PRÉDICTION DES TROUBLES DU CONTRÔLE DE L'IMPULSIVITÉ DANS LA MALADIE DE PARKINSON

par JOHANN FAOUZI

Les troubles du contrôle de l'impulsivité sont une classe de troubles psychiatriques caractérisés par des difficultés dans la maîtrise de ses émotions, pensées et comportements. Ces troubles sont courants dans la maladie de Parkinson et associés à une baisse de la qualité de vie des patients ainsi qu'à une augmentation de la charge des aidants. Pouvoir prédire quels sont les sujets les plus à risque de développer ces troubles et quand ces troubles apparaissent est de grande importance.

L'objectif de cette thèse est d'étudier les troubles du contrôle de l'impulsivité dans la maladie de Parkinson à partir des approches statistique et de l'apprentissage automatique, et se divise en deux parties. La première partie consiste à analyser la performance prédictive de l'ensemble des facteurs associés à ces troubles dans la littérature. La seconde partie consiste à étudier l'association et l'utilité d'autres facteurs, en particulier des données génétiques, pour améliorer la performance prédictive.

Dans un premier chapitre, nous présentons la maladie de Parkinson et les troubles du contrôle de l'impulsivité, effectuons une revue de la littérature sur les troubles du contrôle de l'impulsivité dans la maladie de Parkinson, introduisons les principaux concepts de l'apprentissage automatique et présentons les bases de données sur lesquelles nous avons travaillé et les logiciels utilisés pour analyser ces données. Dans un deuxième chapitre, nous étudions la performance prédictive de plusieurs algorithmes d'apprentissage automatique en utilisant comme variables d'entrée les facteurs associés aux troubles du contrôle de l'impulsivité dans la maladie de Parkinson. Dans un troisième chapitre, nous étudions l'association entre les troubles du contrôle de l'impulsivité dans la maladie de Parkinson et des scores de risque génétique pour un large panel de phénotypes. Dans un dernier chapitre, nous étudions différentes approches d'intégrer des données statiques dans des réseaux de neurones récurrents et évaluons leur performance dans le cas de la prédiction des troubles du contrôle de l'impulsivité dans la maladie de Parkinson, en utilisant des données génétiques pour les données statiques.

À travers ces travaux, nous mettons en avant l'importance d'utiliser des algorithmes d'apprentissage automatique, des méthodes de validation croisée et des cohortes de réplication pour évaluer la puissance prédictive de facteurs de risque connus ou supposés des troubles du contrôle de l'impulsivité dans la maladie de Parkinson.

Remerciements

Je souhaite tout d’abord remercier mes directeurs de thèse, Olivier Colliot et Jean-Christophe Corvol, pour m’avoir donné l’opportunité d’effectuer ce doctorat. Je me trouvais à un moment de ma vie où je n’étais pas sûr de ce que je souhaitais faire, et je leur en serai pour toujours reconnaissant. Je les remercie pour leur encadrement et leur disponibilité, mais aussi pour m’avoir laissé la liberté d’avancer à mon rythme et de travailler sur d’autres projets pendant ma thèse.

Mes remerciements vont également à l’ensemble des membres du jury pour s’être intéressés à mon travail. Je tiens à remercier Marta Avalos et Mélanie Prague pour avoir accepté d’être examitrices, ainsi que Gérard Biau et Alexandre Gramfort pour avoir accepté d’être examinateurs. Je remercie tout particulièrement Christophe Ambroise et Franck Durif pour avoir accepté d’être rapporteurs et pour avoir suivi mon travail tout au long de ma thèse en tant que membres des comités de suivi.

Je remercie les différents membres des équipes cliniques avec qui j’ai travaillé, notamment Stéphanie Carvahlo et Carole Dongmo pour leur aide sur la base de données NS-Park. Je remercie tout particulièrement Samir Bekadar pour toute son aide sur le traitement des bases de données et des données génétiques. Je remercie également Lydia Chougar et Stéphanne Lehericy pour m’avoir permis de travailler sur un projet de classification automatique des syndromes parkinsoniens.

Je remercie bien évidemment tous les membres de l’équipe Aramis qui ont rendu cette thèse fort agréable. Je souhaite remercier Baptiste pour son encadrement et Ninon pour ses conseils. Je remercie notamment mes plus proches voisins : Vincent pour nos discussions matinales, et Tiziana et Giulia pour avoir essayé de m’apprendre l’italien sans grand succès. Je remercie également Elina, Adam et Simona pour les innombrables discussions. Je tiens à remercier en particulier Pascal (le plus gros fly de France) pour ~~m’avoir fait perdre la raison un bon million de fois et des dizaines d’heures à cause du PvP~~ pour son aide et pour les bons moments passés ensemble sur notre passion commune.

Enfin, je souhaite remercier l’ensemble de ma famille, en particulier ma sœur, pour leur soutien indéfectible.

Scientific production

Journal papers

- **Johann Faouzi** and Hicham Janati. “pyts: A Python Package for Time Series Classification”. *Journal of Machine Learning Research*, 21(46):1–6, 2020. <http://www.jmlr.org/papers/v21/19-763.html>
- Romain Tavenard, **Johann Faouzi**, Gilles Vandewiele, Felix Divo, Guillaume Androz, Chester Holtz, Marie Payne, Roman Yurchak, Marc Rußwurm, Kushal Kolar and Eli Woods. “Tslearn, A Machine Learning Toolkit for Time Series Data”. *Journal of Machine Learning Research*, 21(118):1–6, 2020. <http://www.jmlr.org/papers/v21/20-091.html>
- Lydia Chougar, **Johann Faouzi**, Nadya Pyatigorskaya, Rahul Gaurav, Emma Biondetti, Marie Villotte, Romain Valabregue, Jean-Christophe Corvol, Alexis Brice, Louise-Laure Mariani, Florence Cormier, Marie Vidailhet, Gwendoline Dupont, Ines Piot, David Grabli, Christine Payan, Olivier Colliot, Bertrand Degos and Stephane Lehericy. “Automated categorization of parkinsonian syndromes using MRI in a clinical setting”. Accepted in *Movement Disorders*.
- Manon Ansart, Stéphane Epelbaum, Giulia Bassignana, Alexandre Bône, Simona Bottani, Tiziana Cattai, Raphaël Couronné, **Johann Faouzi**, Igor Koval, Maxime Louis, Elina Thibeu-Sutre, Junhao Wen, Adam Wild, Ninon Burgos, Didier Dormont, Olivier Colliot and Stanley Durrleman. “Predicting the Progression of Mild Cognitive Impairment Using Machine Learning : A Systematic and Quantitative Review”. *Medical Image Analysis*, 67(101848), 2021. <https://doi.org/10.1016/j.media.2020.101848>.

Submitted journal papers

- **Johann Faouzi**, Samir Bekadar, Baptiste Couvy-Duchesne, Fanny Artaud, Alexis Elbaz, Graziella Mangone, Olivier Colliot, Jean-Christophe Corvol. “Prediction of impulse control disorders in Parkinson’s disease with replication in an independent cohort”. Submitted in *Annals of Neurology*.
- **Johann Faouzi**, Baptiste Couvy-Duchesne, Samir Bekadar, Olivier Colliot, Jean-Christophe Corvol. “Exploratory analysis of the genetics of impulse control disorders in Parkinson’s disease using genetic risk scores”. Submitted in *Parkinsonism and Related Disorders*.
- Ninon Burgos,* Simona Bottani,* **Johann Faouzi**,* Elina Thibeau-Sutre* and Olivier Colliot. “Deep learning in brain disorders: from data processing to disease treatment”. Submitted in *Briefings in Bioinformatics*.
- Baptiste Couvy-Duchesne,* Johann Faouzi,* Benoît Martin,* Elina Thibeau-Sutre,* Adam Wild,* Manon Ansart, Stanley Durrleman, Didier Dormont, Ninon Burgos, Olivier Colliot. “Ensemble learning of CNN, SVM and BLUP for brain age prediction: ARAMIS contribution to the PAC2019 challenge”. Submitted in *Frontiers in Psychiatry*.

Conference abstracts and posters

- **Johann Faouzi**, Samir Bekadar, Olivier Colliot and Jean-Christophe Corvol. “Predicting Impulse Control Disorders in Parkinson’s Disease: A Challenging Task”. *International Congress of Parkinson’s Disease and Movement Disorders 2019*. <https://hal.inria.fr/hal-02315533v1>

Workshops

- Organization of a one-day workshop entitled [Introduction to machine learning in Python with Scikit-learn](#) at the Paris Brain Institute, as part of the [Partnership for Advanced Computing in Europe](#).
- Organization of a practical session on *Image Synthesis using Generative Adversarial Networks* at the [Hands-on Workshop on Machine Learning Applied to Medical Imaging](#) taking place at the Paris Brain Institute..

Teaching

- Teaching assistant for the *Python Programming for mathematics* course at Sorbonne Université, France ([Programmation Python pour les mathématiques](#), in French), 2019-2020.
- Teaching assistant for the [Deep Learning for Medical Imaging](#) course at Master MVA, France, 2019-2020.

Contents

Abstract	iii
Résumé	v
Remerciements	vii
Scientific production	ix
List of Figures	xvii
List of Tables	xix
List of Abbreviations	xxi
Introduction	1
1 Background	5
1.1 Parkinson's disease	6
1.1.1 History	6
1.1.2 Classification	7
1.1.3 Pathophysiology	7
1.1.4 Diagnosis	8
1.1.5 Symptoms	9
1.1.6 Medications and their limitations	11
1.1.7 Motor complications	15
1.2 Impulse control disorders	16
1.2.1 Definition of specific impulse control disorders	16
1.2.2 Studies on impulse control disorders in subpopulations	19
1.3 Impulse control disorders in Parkinson's disease	19
1.3.1 Epidemiology	20
1.3.2 Assessment and diagnosis	20
1.3.3 Associations	21
1.3.4 Prediction	23

1.3.5	Other behavioral addictions	24
1.4	Machine learning	25
1.4.1	Notations	25
1.4.2	Algorithms	26
1.4.3	Regularization	36
1.4.4	Metrics	42
1.5	Putting it all together	46
1.6	Materials	47
1.6.1	Data sets	48
1.6.2	Software	50
2	Prediction of impulse control disorders in Parkinson's disease	55
2.1	Introduction	56
2.2	Materials and methods	57
2.2.1	Populations	57
2.2.2	Participants and clinical measurements	57
2.2.3	Genetic variants	58
2.2.4	Data processing	58
2.2.5	Machine learning algorithms	58
2.2.6	Cross-validation	60
2.2.7	Statistical analysis	60
2.3	Results	61
2.3.1	Population characteristics	61
2.3.2	Predictive performance	62
2.3.3	Contribution of the different features	64
2.4	Discussion	70
3	Exploratory analysis of the genetics of impulse control disorders in Parkinson's disease using genetic risk scores	73
3.1	Introduction	74
3.2	Materials and methods	75
3.2.1	Populations	75
3.2.2	Participants	75
3.2.3	Genetic ancestry	76
3.2.4	Genotyping and quality control	76
3.2.5	Phenotypes and genome-wide association studies	77
3.2.6	Computation of genetic risk scores	77
3.2.7	Statistical analyses	77
3.3	Results	78
3.3.1	Participants and genetic variants	78

3.3.2	Genome-wide association studies	78
3.3.3	Association analyses	78
3.4	Discussion	82
4	Combining static and dynamic data in recurrent neural networks	85
4.1	Introduction	85
4.2	Related work	86
4.3	Proposed approach	89
4.4	Experiments	91
4.5	Conclusion	93
	Conclusion	95
A	Supplementary materials for the prediction of impulse control disorders from clinical and genetic data with replication in an independent cohort	99
A.1	Reduction approaches	99
A.2	Supplementary Tables	100

List of Figures

1.1	Schematic overview of the primary motor circuits in the basal ganglia, the indirect (left) and direct (right) pathways	8
1.2	Evidence-based medicine review of treatment options for motor symptoms of early PD	12
1.3	Evidence-based medicine review of treatment options for motor symptoms of treated PD optimized on levodopa	13
1.4	Ordinary least squares linear regression	27
1.5	Decision function of a logistic regression model	28
1.6	Support vector machine classifier	29
1.7	Impact of the kernel on the decision function of a support vector machine classifier	30
1.8	A decision tree	31
1.9	Example of an artificial neural network	33
1.10	Main concept of a recurrent neural network	34
1.11	Vanilla recurrent neural network unit	35
1.12	Long Short-Term Memory unit	36
1.13	Gated Recurrent Unit	37
1.14	Relationship between error and model complexity	38
1.15	Toy regression data set with a non-linear relationship	38
1.16	Illustration of regularization	40
1.17	Unit balls of the ℓ_0 , ℓ_1 and ℓ_2 norms	41
1.18	Receiver operating characteristic and precision-recall curves	46
2.1	Architecture of the recurrent neural network	59
2.2	Cross-validation procedure	60
2.3	ROC and Precision-recall curves on PPMI	65
2.4	ROC and Precision-recall curves on DIGPD	66
2.5	Statistical comparison of ROC AUC for the four main models	67
3.1	Genetic ancestry estimation	79

4.1	Recurrent neural network with no static data	87
4.2	Recurrent neural network with static data on its own branch	88
4.3	Recurrent neural network with static data treated as dynamic data . . .	88
4.4	Recurrent neural network with static data initializing the GRU layer . .	89
4.5	Recurrent neural network with static data modifying the dynamic features	90
4.6	Cross-validation procedure	93

List of Tables

1.1	Confusion matrix for binary classification	43
1.2	PPMI clinical sites	49
1.3	Questionnaires and scales used in PPMI and DIGPD	51
2.1	Baseline characteristics	62
2.2	Results of the four main models	63
2.3	Statistical comparison of ROC AUC for the four main models with and without genetic variants	68
2.4	Coefficients of the three logistic regression models without genetic variants as input	69
3.2	Results of the association analyses	81
4.1	Genome-wide association studies from which genetic risk scores were derived	92
4.2	Predictive performance of the five approaches	93
A.1	Genetic variants included in the analyses	101
A.2	Coefficients of the three logistic regression models with genetic variants as input	102
A.3	Predictive performance on PPMI of the five machine learning algorithms with the nine reduction approaches	103
A.4	Predictive performance on DIGPD of the five machine learning algorithms with the nine reduction approaches	104
A.5	Predictive performance on PPMI of the five machine learning algorithms with the nine reduction approaches with 10 repetitions of the nest cross-validation	105
A.6	Predictive performance on DIGPD of the five machine learning algorithms with the nine reduction approaches with 10 repetitions of the nest cross-validation	106

List of Abbreviations

AP	Average precision
COMT	Catechol-O-methyltransferase
DA	Dopamine agonist
DIGPD	Drug Interaction With Genes in Parkinson’s Disease
FN	False negatives
FP	False positives
GRS	Genetic risk score
GRU	Gated Recurrent Unit
GWAS	Genome-wide association study
ICD	Impulse control disorder
LDR	Long duration release
LSTM	Long Short-Term Memory
MAO	Monoamine oxidase
MDS-UPDRS	Movement Disorder Society-sponsored revision of the Unified Parkinson’s Disease - Rating Scale
NPV	Negative predictive value
PD	Parkinson’s Disease
PPMI	Parkinson’s Progression Markers Initiative
PPV	Positive predictive value
PR	Precision-recall
QUIP	Questionnaire for Impulsive-Compulsive Disorders in Parkinson’s Disease
QUIP-RS	Questionnaire for Impulsive-Compulsive Disorders in Parkinson’s Disease - Rating Scale
RNN	Recurrent neural network
ROC	Receiver operating characteristic
ROC AUC	Area under the receiver operating characteristic curve
SDR	Short duration release
SNP	Single nucleotide polymorphism
SVM	Support vector machine
TN	True negatives
TNR	True negative rate
TP	True positives
TPR	True positive rate
UMAP	Uniform Manifold Approximation and Projection

Introduction

Context

Parkinson's disease (PD) is a neurodegenerative disease with no cure to date. Besides its characteristic motor symptoms, numerous non-motor symptoms have been reported to occur in the course of the disease ([Chaudhuri et al., 2006](#)). A specific symptom has recently been increasingly acknowledged: impulse control disorders.

Impulse control disorders (ICDs) are a class of psychiatric disorders involving problems in the self-control of emotions and behaviors ([American Psychiatric Association, 2013](#)). They include pyromania, kleptomania, Internet addiction disorder and intermittent explosive disorder for instance. In Parkinson's disease, the four most common impulse control disorders are pathological gambling, binge eating disorder, compulsive sexual behavior and compulsive buying disorder ([Weintraub and Claassen, 2017](#)).

Many factors have been associated with impulse control disorders in Parkinson's disease ([Grall-Bronnec et al., 2018](#)). Socio-demographic factors (age of PD onset, sex), psychiatric comorbidities (anxiety, depression), sleep disorders, and PD medication have been associated with ICDs in PD among others. Dopamine agonists (DAs), a class of PD medication, have been the most strongly correlated with ICDs. Dopamine agonists still have their advantages as they delay the initiation of levodopa, another class of PD medication with its own adverse effects.

Impulse control disorders are associated with a decrease in quality of life, strained interpersonal relationships, financial distress, medical complications, and higher caregiver burden ([Weintraub and Claassen, 2017](#)). Prompt identification and treatment of the symptoms are usually imperative to improve the quality of life of the subjects. However, managing impulse control disorders implies that they are already present. Accurately predicting which subjects are at higher risk of developing these disorders and when they occur could allow for early management, decrease their negative impact, and potentially *prevent* their apparition.

Objective

The general objective of this thesis is to study impulse control disorders in Parkinson’s disease from the machine learning point of view, with a focus on their predictability. Despite the rise of machine learning, this approach has been little to not used in the context of ICDs in PD.

Our first objective is to study the how well ICDs in PD can be predicted by combining the factors reported in the literature. The known risk factors come from univariate association studies that do not take into account the other risk factors. Machine learning allows for learning from a set of features at once and thus leveraging information from all the factors. Using machine learning comes with specific methodological requirements to assess their predictive performance in an unbiased manner. These methods, such as cross-validation, are not always well-known in the medical field but are necessary not to report overly optimistic results. We describe and apply such methods in this application.

Our second objective is to investigate the genetic factors of ICDs in PD, as little is known about these risks. A few associations from candidate genes analyses have been reported, but these studies have not been replicated. On the other hand, genome-wide association studies (GWAS) investigate the combined risk of the whole genome by computing a genetic risk score (GRS), but none has been published on ICDs in PD. Instead of performing a GWAS, which requires a large sample size, we study the association between ICDs in PD and genetic risk scores of other phenotypes, for which large GWAS exist.

Contributions

Our contributions in the field of impulse control disorders in Parkinson’s disease are three-fold. First, we study the added value of combining the reported factors associated with ICDs in PD to predict them, by training machine learning algorithms using these factors as input. Second, we investigate the association between ICDs in PD and genetic factors for a broad range of phenotypes, including other psychiatric disorders. Third, we study how to integrate time-dependent features, such as clinical measurements, and time-independent information, such as socio-demographic and genetic factors, in predictive models, with an application in the prediction of ICDs in PD using recurrent neural networks.

Only two studies have reported a classification task of impulse control disorders in Parkinson’s disease (Erga et al., 2018; Kraemmer et al., 2016). However, both studies lack a replication cohort and have major methodological issues (lack of cross-validation, biased feature selection) that alter the trust in the reported predictive performance. We propose the first study evaluating in an unbiased manner the predictive performance

of several machine learning algorithms using known risk factors as input data. We investigate the use of five standard machine learning classification algorithms (logistic regression, support vector machines with linear and RBF kernels, random forest and gradient tree boosting) and recurrent neural networks to predict the presence or absence of ICDs for a given patient at their next visit. In order to make the variable-length sequences of visits suitable as input data for the standard machine learning algorithms, we reduce each variable-length sequence of visits into one “summary” visit using a convex combination. We investigate several reduction approaches, each giving different weights to each visit. We evaluate the predictive performance on two research cohorts with different characteristics to assess the generalization capability of the models.

Secondly, we investigate the association between genetic risk scores and impulse control disorders in Parkinson’s disease. Many phenotypes are known to be heritable, yet less is known about which parts of the genomes and how they contribute to this heritability. A genetic risk score is a single score indicating, given one’s genome, their risk of developing a given phenotype (Wray et al., 2007). For instance, genetic risk scores for schizophrenia and bipolar disorder have been associated with creativity (Power et al., 2015). Recently, the genetic risk score of Parkinson’s disease has been reported not to be associated with ICDs in PD (Ihle et al., 2020). We investigate the association between 40 generic risk scores representing a broad range of phenotypes, and ICDs in PD, in two research cohorts.

Finally, we investigate the integration of static data in recurrent neural networks. We review the literature on this topic and identify four approaches. Besides the dummy approach consisting in removing static data, the most common approach is to have static and dynamic data in their own branches in the network. This approach is not specific to combining static and dynamic data, and is commonly used to integrate multimodal data in artificial neural networks (Hao et al., 2019; Mobadersany et al., 2018). The two other identified approaches consist in treating static data as dynamic data (Leontjeva and Kuzovkin, 2016; Rahman et al., 2020) and initializing the parameters of the recurrent neural unit using the static features (Kristensen and Burelli, 2019). We propose a new approach consisting in modifying the dynamic features using the static features. We investigate the predictive performance of the five approaches in the use case of predicting ICDs in PD, where the static features are socio-demographic and genetic features.

Chapter 1

Background

Parkinson's disease is the second most frequent neurodegenerative disorder after Alzheimer's disease (Nussbaum and Ellis, 2003). In France, 150,000 people were affected by PD in 2010 and over 250,000 are expected to be affected in 2030 (Wanneveich et al., 2018). To date, no cure for this disease exists and the quality of life of the patients slowly but steadily decreases until death. Its economic impact is substantial as the cost in the United Kingdom was estimated to be between £449 million and £3.3 billion annually in 2007 (Findley, 2007). The social impact of PD is also important as the quality of life of PD cases is heavily reduced due to the large range of impairing symptoms, increasing caregiver burden.

Besides the three cardinal motor symptoms of PD that are tremor, bradykinesia, and rigidity (Jankovic, 2008), many non-motor symptoms often occur during the course of the disease including cognitive, sleep, dysautonomic, and behavioral disorders (Chaudhuri et al., 2006). There is no cure for PD. The dopamine replacement therapy alleviate motor symptoms, but is hampered by motor complications (fluctuations, dyskinesia) and adverse effects. Impulse control disorders, a class of behavioral psychiatric disorders characterized by impulsivity, is a frequent adverse effect of replacement dopamine therapy (Corvol et al., 2018). Predicting which subjects will develop these disorders and when is challenging, but of great importance because of their familial, social, economic or legal impact.

In this section, we introduce Parkinson's disease, from pathophysiology to symptoms to medications. Then we define impulse control disorders and detail in which populations they have been studied, and review the literature on impulse control disorders in Parkinson's disease. Next we introduce the main concepts of machine learning. Finally we discuss how machine learning can be helpful to tackle impulse control disorders in Parkinson's disease.

1.1 Parkinson's disease

1.1.1 History

The first clear clinical description of Parkinson's disease was provided in 1817 by James Parkinson in his *Essay on the Shaking Palsy* (Parkinson, 2002), in which he defined shaking palsy as:

Involuntary tremulous motion, with lessened muscular power, in parts not in action and even when supported; with a propensity to bend the trunk forwards, and to pass from a walking to a running pace: the senses and intellects being uninjured

and used the term *paralysis agitans* to describe individuals with this disorder. In his essay, he reported six cases that he had seen as patients or that he had observed during his wanderings through the streets near his home in Hoxton Square (Lees, 2007).

Jean-Martin Charcot, the father of modern neurology (Gomes and Engelhardt, 2013), deepened the knowledge on this disorder with his studies between 1868 and 1881, notably distinguishing between rigidity, weakness and bradykinesia (Lees, 2007). He also advocated the renaming of the disorder in honor of James Parkinson (Lees, 2007).

Landmarks on the understanding of the disease include the description of microscopic particles, later called Lewy bodies, in the brains of PD cases by Friedrich Lewy in 1912 (Holdorff, 2006), the report that the substantia nigra was the main structure affected in the brain by Konstantin Tretiakoff in 1919 (Goedert et al., 2013), the underlying biochemical changes in the brain by Arvid Carlsson and Oleh Hornykiewicz in the 1950s (Lees, 2007), and the discovery of alpha-synuclein being the main component of Lewy bodies by Spillantini and others in 1997 (Spillantini et al., 1997).

Antiparkinsonian effects of anticholinergics were first described in 1868 by Leopold Ordenstein, a student of Jean-Martin Charcot (Kim et al., 2017). In his thesis, he states that "Monsieur Charcot has begun to prescribe 2 or 3 granules of hyoscyamine daily, approximately 1 mg each. This medication was able to provide several hours of rest for some patients. Apparently, further observations are necessary to make a decision about this medication". In 1957, Arvid Carlsson showed that levodopa reversed the akinetic effect of reserpine, a drug that lowers blood pressure and slows heart rate, in rabbits (Carlsson et al., 1957). In 1960, Oleh Hornykiewicz published a landmark paper showing for the first time a significant depletion of dopamine in the caudate and putamen of patients only with PD or postencephalitic parkinsonism (Ehringer and Hornykiewicz, 1960). The successful introduction of high dosage levodopa therapy occurred in 1967 (Cotzias et al., 1967). Before, anticholinergics remained the only available medical therapy for Parkinson's disease. The 1960s were also marked by the first observation of the antiparkinsonian effects of amantadine. In the 1980s, dopamine agonists were tested

as monotherapy in early PD and two catechol-O-methyltransferase (COMT) inhibitors were found to be orally active, following the discovery of monoamine oxidase (MAO) as the mechanism for inactivating the monoamines levodopa decarboxylase enzyme in the 1930s (Kim et al., 2017).

1.1.2 Classification

In the tenth revision of the International Statistical Classification of Diseases and Related Health Problems, the code for Parkinson's disease is G20, belonging to the group of extrapyramidal and movement disorders (G20–G26) among the diseases of the nervous system (G00–G99).

The main motor symptoms of PD, called parkinsonism, consist of bradykinesia and one of two other physical signs: muscular rigidity and tremor at rest (Jankovic, 2008). Parkinson's disease is the most prevalent form of parkinsonism and is often called *idiopathic parkinsonism*, that is parkinsonism with no identifiable cause (Samii et al., 2004). Drugs, toxins, infections, and brain lesions such as stroke can lead to parkinsonism. Parkinsonism is not exclusive to PD and can be found in a group of other diseases often called atypical parkinsonism, consisting of other features differentiating them from PD. This group includes multiple system atrophy, progressive supranuclear palsy, Lewy body dementia and corticobasal degeneration (Samii et al., 2004).

1.1.3 Pathophysiology

The pathological hallmark of PD is cell death in the basal ganglia, particularly in the ventral component of the substantia nigra pars compacta (Davie, 2008). By the time of death, the substantia nigra pars compacta has lost up to 70% of its neurons in comparison to unaffected individuals. Death of astrocytes and a significant increase in the number of microglia in the substantia nigra also occur (Dickson, 2018).

Figure 1.1 illustrates the primary motor circuits in the basal ganglia. The basal ganglia are functionally connected to other brain regions via the motor, oculo-motor, associative, limbic, and orbitofrontal pathways (Obeso et al., 2008). The motor pathway connects the basal ganglia to the motor cortex, which is involved in the planning, control, and execution of voluntary movements. The oculo-motor pathway links the basal ganglia to the frontal eye fields, which is responsible for saccadic and voluntary eye movements. The cerebral cortex is connected to the basal ganglia via the associative pathway, and enable to support abstract thinking and language, produce a meaningful perceptual experience of the world, and enable us to interact effectively. The limbic pathway connects the basal ganglia to the limbic system, which supports a variety of functions including emotion, behavior, motivation, long-term memory, and olfaction. The orbitofrontal cortex is connected to the basal ganglia via the orbitofrontal pathway, which is involved in the cognitive process of decision-making.

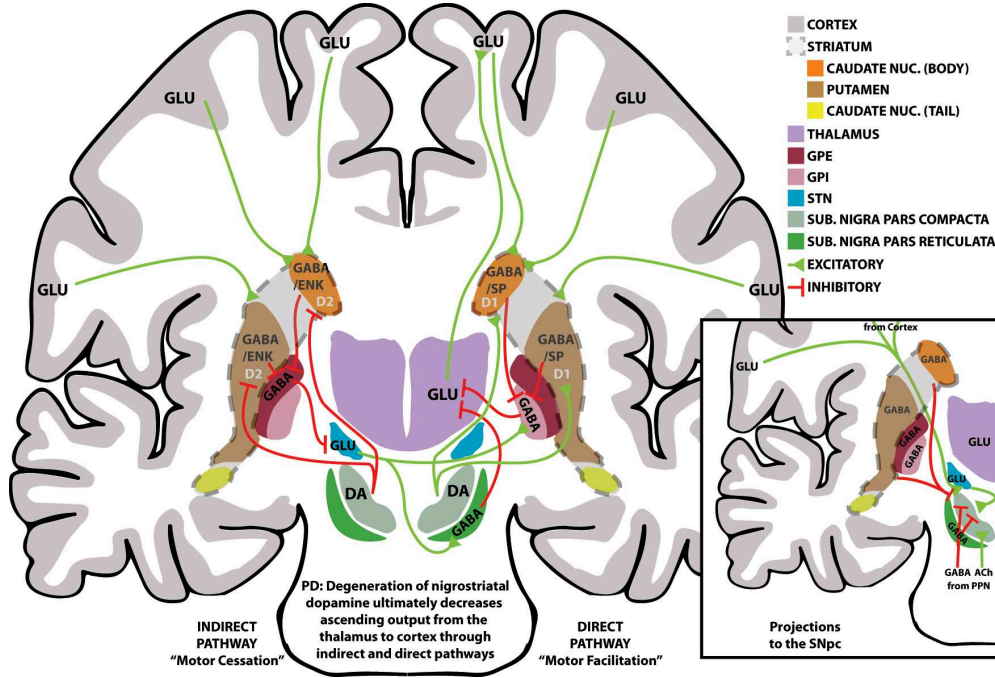


Figure 1.1: Schematic overview of the primary motor circuits in the basal ganglia, the indirect (left) and direct (right) pathways (reproduced from (Harris et al., 2020b)).

D1: D1 receptors; D2: D2 receptors; DA: dopamine; ENK: enkephalin; GLU: glutamate; NUC: nucleus; PPN: pedunculopontine nucleus; SP: substance P; SUB: substantia.

Dopamine neurons of the substantia nigra pars compacta mainly innervate the motor pathway, hence the predominance of the motor symptoms. However, the mesolimbic and mesocortical pathways have their somas in the ventral tegmental area, which is relatively sparsely impacted in PD. The motor symptoms are actually assumed to be caused by the impairment of these pathways rather than the motor pathway (Obeso et al., 2008).

1.1.4 Diagnosis

A definite diagnosis of Parkinson's disease requires autopsy, the final proof being the presence of Lewy bodies in the midbrain (Samii et al., 2004). Nonetheless, clinical diagnosis of this disorder has become more rigorous and several criteria have been proposed (Calne et al., 1992; Hughes et al., 1992; Postuma et al., 2015; Ward and Gibb, 1990). They include confidence levels such as clinically possible, clinically probable and clinically definite, and rely on the presence of parkinsonism, response to antiparkinson drugs and exclusion criteria that would favor another cause for the presence of parkinsonism (Samii et al., 2004).

1.1.5 Symptoms

The symptoms of Parkinson's disease are numerous and have a substantial negative impact on the quality of life of the individuals (Jankovic, 2008). Apart from the recognizable motor symptoms, many other symptoms occur frequently during the course of the disease. Particularly, a wide range of neuropsychiatric disorders have been reported in individuals with PD (Balestrino and Martinez-Martin, 2017).

Motor symptoms

The three cardinal motor symptoms of Parkinson's disease are bradykinesia, rigidity, and rest tremor (Jankovic, 2008). Other typical symptoms of parkinsonism include postural instability, flexed posture and freezing usually occurring later in the course of the disease and leading to falls (Jankovic, 2008; Kalia and Lang, 2015).

Akinesia, bradykinesia and hypokinesia are the hallmarks of basal ganglia disorders, and includes difficulties with planning, initiating and executing movement, and with performing sequential and simultaneous tasks. Bradykinesia refers to slowness of movement and is the most characteristic clinical feature of PD. Bradykinesia may be associated with hypokinesia (reduction in movement amplitude) or akinesia (poverty of action and difficulty initiating movements) (Moustafa et al., 2016).

The initial manifestation is usually slowness in performing activities of daily living, slow movement and reaction times (slowness of walking), including difficulties with tasks requiring fine motor control such as handwriting (micrographia), using utensils and buttoning. Assessment of bradykinesia usually includes having patients perform rapid, repetitive, alternating movements of the hand and heel taps and observing slowness and decreasing amplitude (Jankovic, 2008).

Rigidity is characterised by increased resistance and is usually accompanied by the "cogwheel" phenomenon, that is a circular jerking rigidity in flexion and extension in a background of tremor. Rigidity can occur at many locations, including neck, shoulders, hips, wrists and ankles. Reinforcing manoeuvres tend to increase rigidity and can be used to detect mild cases of rigidity (Jankovic, 2008).

Tremor at rest is the most recognizable symptom of PD. Tremors are typically unilateral at disease onset, occur in the rest position at a frequency between 4 and 6 Hz, disappear during action, and are usually prominent in the distal part of extremities (Jankovic, 2008). In particular, hand tremors consist of the tendency of the thumb and the index finger to approximate one another while trembling as if an object was being rolled between the two fingers. The term "pill-rolling" is often used to describe these tremors because of the similarity to the technique used by pharmacists to fashion a pill by rolling a substance between the two fingers (Cooper et al., 2008). Rest tremor in patients with PD can also involve the lips, chin, jaw and legs (Jankovic, 2008).

Postural instability due to loss of postural reflexes is usually a symptom of the late

stages of PD and generally occurs after the onset of other clinical features. The pull test is often used to assess postural instability: the patient is quickly pulled backward by the shoulders to assess the degree of retropulsion. Taking more than two steps backwards or the absence of any postural response indicates an abnormal postural response. Postural instability and freezing of gait are the most common causes of falls and contribute substantially to the risk of hip fractures ([Jankovic, 2008](#)).

Rigidity and postural deformities can result in flexed posture (camptocormia), such as flexed neck, and usually occur in a late stage of the disease ([Jankovic, 2008](#)).

Although the dopamine denervation is always bilateral in PD, it is commonly asymmetric, and the symptoms typically occur unilaterally at the onset of the disease, affecting the other part of the body during disease progression. Akinesia, tremor, and rigidity are responsive to the dopamine replacement therapy which typically improve the symptoms by more than 70% in most patients. By contrast, postural instability and gait disturbance are not poorly responsive to treatment.

Non-motor symptoms

Many non-motor symptoms have also been reported in Parkinson's disease and can impact the quality of life of the subjects and caregiver burden more than the motor symptoms ([Hiseman and Fackrell, 2017](#)). Indeed, these symptoms are not responsive to the treatment and may even worsen with dopamine replacement therapy, and largely contribute to the burden of the disease for the patients and the caregivers. Neuropsychiatric symptoms are particularly common during the course of the disease ([Balestrino and Martinez-Martin, 2017](#)).

Major depressive disorder is frequent in PD, with an approximated 17% prevalence. Comorbid depression worsens cognition, function, and quality of life, and increases caregiver burden and mortality. Symptomatic overlap between major depression disorder and PD can make appropriate detection and treatment difficult ([Goodarzi et al., 2016](#)).

Up to 55% of PD patients experience substantial anxiety symptoms, and up to 40% have an anxiety disorder as defined by the criteria of the Diagnostic and Statistical Manual of Mental Disorders. The most common anxiety disorders in PD are generalized anxiety disorder, and social and other phobias ([Broen et al., 2016](#)).

Most PD patients suffer from cognitive decline or dementia during the course of the disease. The prevalence of mild cognitive impairment is around 25% in individuals with PD but without dementia ([Litvan et al., 2011](#)). The point prevalence of PD dementia is approximately 30% and its cumulative prevalence is at least 75% for PD patients surviving more than 10 years ([Litvan et al., 2011](#)). Cognitive impairment mostly affect executive and visuo-spatial functions, rather than memory disturbances, and heavily impacts functioning, caregiver burden and mortality ([Goldman et al., 2018](#)).

Sleep disturbances are common in PD and consist mainly of nighttime sleep difficulties such as insomnia, restless legs syndrome, rapid eye movement sleep behavior

disorder, and sleep-disordered breathing, but also of excessive daytime sleepiness. The prevalence of insomnia, based on physician interview, is estimated to be around 30–59% (Chahine et al., 2017). Daytime sleepiness can make individuals with PD quit driving, increasing caregiver burden. Sleepiness can be related to the disease itself, but is also an adverse effect of dopamine replacement therapy.

The majority of PD subjects also suffer from gastrointestinal symptoms, constipation being considered the most prominent (Mertsalmi et al., 2017). Other gastrointestinal symptoms include drooling, taste impairment, swallowing disorders (Fasano et al., 2015), and irritable bowel syndrome (Mertsalmi et al., 2017).

Apart from constipation, other autonomic dysfunction occurs frequently, the most common symptoms being orthostatic hypotension, urinary and sexual dysfunction, abnormal sweating and seborrhoea (Jankovic, 2008).

A number of neuro-ophthalmological abnormalities may be seen in patients with PD, including visual hallucinations, ocular surface irritation, decreased blink rate, altered tear film, blepharospasm and decreased convergence (Biouesse et al., 2004).

Impulse control and related behaviors are common comorbidities and are strongly associated with dopamine replacement therapy. Almost half of PD patients are expected to develop impulse control disorders five years after PD onset (Corvol et al., 2018). The four major ICDs that have been reported in PD are pathological gambling, compulsive shopping, binge eating and hypersexuality (Weintraub and Claassen, 2017). Other related impulsive-compulsive behaviors include dopamine dysregulation syndrome (Cilia et al., 2014), punding (Evans et al., 2004) and hobbyism (Callesen and Damholdt, 2017).

1.1.6 Medications and their limitations

Contrary to Alzheimer's disease, for which there exists no treatment that substantially decreases the magnitude of the main symptoms, several therapies are effective at limiting the decrease in quality of life of individuals with PD (Fahn, 2008). The most simple yet efficient therapy is dopamine replacement therapy, which consists in replacing the loss of dopamine due to the cell death in the basal ganglia. Its main classes of medications consist of levodopa, dopamine agonists, and inhibitors. Other therapies include deep brain stimulation (Herrington et al., 2016) and exercise programs (Ahlskog, 2011). **Figure 1.2** and **Figure 1.3** summarize the treatment options for PD.

Levodopa Levodopa is an abbreviation of L-3,4-dihydroxyphenylalanine and is the precursor to dopamine (Fahn, 2008). Contrary to levodopa, dopamine itself is unable to cross the blood-brain barrier and cannot be used to treat PD (Zahoor et al., 2018). After absorption and transit across the blood-brain barrier, levodopa is converted into the neurotransmitter dopamine by DOPA decarboxylase. Patients are usually administered low dose of levodopa, with the dose being adjusted based on the patient's response to treatment and balanced against the adverse effects experienced. Although levodopa

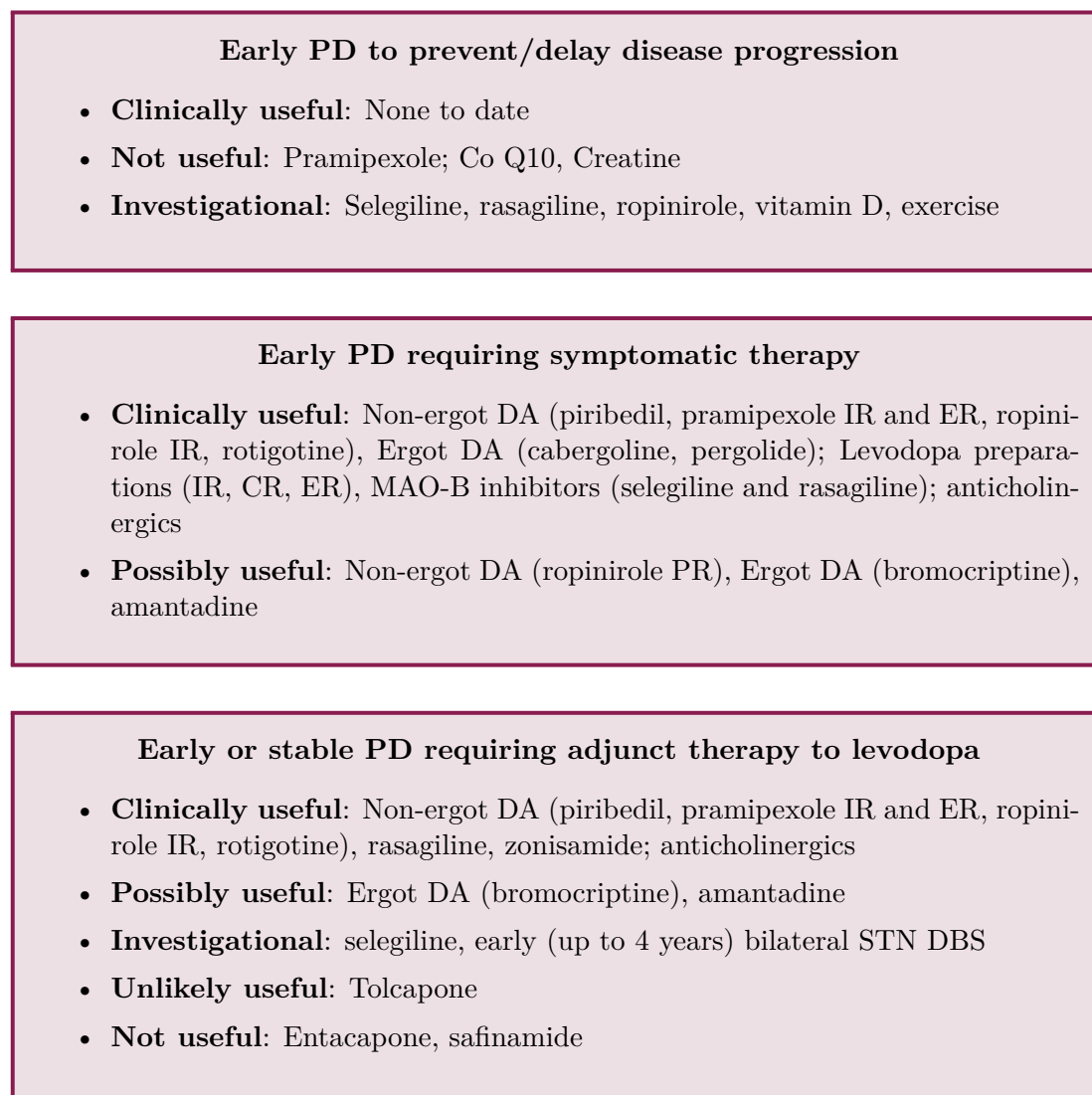


Figure 1.2: Evidence-based medicine review of treatment options for motor symptoms of early PD (reproduced from (Fox et al., 2018)).

CR: controlled release; DA: dopamine agonist; DBS: deep brain stimulation; ER: extended release; IR: immediate release; PR: prolonged release; STN: subthalamic nucleus.

Treating motor fluctuations

- **Clinically useful:** Non-ergot DA (pramipexole, ropinirole, rotigotine, apomirphine intermittent injections, pergolide); levodopa ER; COMT inhibitors (entacapone; opicapone); MAO-B inhibitors (rasagiline, safinamide, zonisamide); LCIG; bilateral DBS (STN or GPi)
- **Possibly useful:** Ergot DA (bromocriptine, cabergoline); istradefylline; tolcapone; Non-ergot DA (apomorphine infusion)
- **Investigational:** Selegiline, rasagiline, ropinirole, vitamin D, exercise

Treating dyskinesia

- **Clinically useful:** Amantadine; clozapine; LCIG, bilateral DBS surgery (STN or GPi); unilateral pallidotomy

Treating specific/general motor symptoms

- **Clinically useful:** Physiotherapy
- **Possibly useful:** Rivastigmine (gait and balance); Exercise-based movement strategy training (gait and balance); formalized patterned exercises (gait and balance); speech therapy (speech and swallowing); occupational therapy; thalamic surgery (DBS or thalamotomy) (tremor)
- **Investigational:** Donepezil (gait and balance); methylphenidate (gait and balance); memantine (gait and balance); cannibidiol; technology-based movement strategies; acupuncture; rTMS; tDCS

Figure 1.3: Evidence-based medicine review of treatment options for motor symptoms of treated PD optimized on levodopa (reproduced from (Fox et al., 2018)).

COMT: catechol-O-methyltransferase; CR: controlled release; DA: dopamine agonist; DBS: deep brain stimulation; ER: extended release; GPi: globus pallidus interna; IR: immediate release; LCIG: levodopa-carbidopa intestinal gel; MAO-B: monoamine oxidase B PR: prolonged release; rTMS: repetitive transcranial magnetic stimulation; STN: subthalamic nucleus; tDCS: tDirect Current Stimulation.

is effective against the main motor symptoms, it comes up with side effects such as dyskinesia (involuntary movements) and fluctuations in effectiveness ([National Clinical Guideline for Diagnosis and Management in Primary and Secondary Care, 2006](#)). Levodopa is also ineffective against several motor (gait, speech) and non-motor (cognitive, sensory, vegetative) PD symptoms ([You et al., 2018](#)).

Dopamine agonists An agonist is a chemical that binds to a receptor and activates the receptor to produce a biological response. Dopamine agonists stimulate the activity of the dopamine system by binding to the dopaminergic receptors and, unlike levodopa, do not need to be converted into dopamine ([Zahoor et al., 2018](#)). Initial treatment with dopamine agonists allows for a delay in the use of levodopa, which may curtail the impact of the problematic motor complications ([Rascol et al., 2000](#)). However dopamine agonists are less potent than levodopa, and are less tolerated than levodopa with higher rate of nausea and vomiting, insomnia, sleepiness, and hallucinations, particularly in the elderly. Ergot DAs have also been progressively abandoned because of their association with heart valve fibrosis ([Corvol et al., 2007](#); [Van Camp et al., 2004](#)). One of the most troublesome characteristic side effects of dopamine agonists is probably impulse control disorders and related behaviors, including pathological gambling, hypersexuality, binge eating, compulsive shopping, punding, and hobbyism.

Inhibitors of levodopa metabolism An inhibitor is a substance that decreases the rate of, or prevents, a chemical reaction. Inhibiting enzymes that are involved in dopamine degradation is the main feature of this class of medications. Monoamine Oxidase B (MAO-B) is one of the main enzymes involved in the breakdown of dopamine, thus reducing the activity of this enzyme results in increased dopaminergic activity within the striatum ([Zahoor et al., 2018](#)). The intake of MAO-B inhibitors relieves motor symptoms in PD patients, and as with dopamine agonists they may be used as an initial treatment option. The main side effects of MAO-B inhibitors are increased dyskinesia and headaches ([Connolly and Lang, 2014](#)). Catechol-O-methyl transferase is another enzyme involved in dopamine as well as in levodopa degradation. COMT inhibitors are used as adjunctive therapy to levodopa by prolonging its duration of action ([Zahoor et al., 2018](#)). The most common adverse effects of both entacapone and tolcapone, the most-used COMT inhibitors, are increased dyskinesia and diarrhoea in up to 20% of the treated patients ([You et al., 2018](#)).

Others Levodopa, dopamine agonists and inhibitors are all designed to increase dopaminergic activity in the striatum. A few drugs that act through non-dopaminergic mechanisms are also used in the treatment of PD ([Zahoor et al., 2018](#)). Anticholinergics, by acting as antagonists at cholinergic receptors, limit the activity of the neurotransmitter acetylcholine. Their most common adverse effects include hallucinations, blurred

vision, dry mouth, constipation, drowsiness, memory problems, and increased dyskinesia (Zahoor et al., 2018). Amantadine, which was initially developed as an antiviral drug for treating flu, has subsequently been used for the treatment of PD. It may be used to treat rigidity, rest tremor, and is also used to treat levodopa-induced dyskinesia (Ory-Magne et al., 2014; Zahoor et al., 2018). While generally well tolerated, possible adverse effects associated with the use of amantadine include hallucinations, confusion, blurred vision, impaired concentration, nausea and vomiting (Zahoor et al., 2018).

1.1.7 Motor complications

Despite its spectacular effects on the core motor symptoms of Parkinson's disease, levodopa is not a perfect drug as it does not fulfil all the needs of PD patients. The long-term outlook for PD patients is hampered by the occurrence of motor complications: motor fluctuations and dyskinesia. Levodopa and inhibitors of levodopa metabolism are notably associated with increased dyskinesia (You et al., 2018). Dyskinesia and motor fluctuations affect virtually all patients but the delay in their occurrence is highly variable. More than 90% of PD patients are expected to experience motor complications after 10 years (Hely et al., 1999; Mazzella et al., 2005).

Motor fluctuations Motor fluctuations are characterised by wearing-off, that is worsening or reappearance of motor symptoms before the next levodopa dose resulting in an “off” state that improves when the next dose is taken (“on” state) (You et al., 2018). There are two kinds of response to levodopa: the short duration response (SDR) and the long-duration response (LDR). The former corresponds to the motor improvement following a single dose of levodopa and lasts from minutes to hours (Muentner and Tyce, 1971). Its effect is immediately lost if levodopa is stopped. The latter has a slower development and builds up during repeated levodopa dosing, taking days to weeks to come into effect, but also decays gradually over a similar span of time after levodopa has been withdrawn (Anderson and Nutt, 2011). Both mechanisms are present from the beginning of PD treatment. The SDR accounts for a half to two thirds of the motor response, while the LDR accounts for the remaining part (Ogasahara et al., 1984). However, their effects are not strictly additive but overlapping and even show a different time course as the disease progresses. In the early stages of PD, the LDR predominates and masks most of the SDR, thus patients have a stable response to levodopa (Nutt and Holford, 1996). With disease progression and long-term levodopa treatment, the LDR decreases and the SDR shortens, with a more immediate onset and decline and a greater difference between baseline and peak response. Therefore, the masking of the SDR by the LDR dwindles and patients experience motor fluctuations (Nutt and Holford, 1996).

Levodopa-induced dyskinesia Dyskinesias are abnormal movements of the limbs, the trunk and the face induced by the dose of levodopa. Many patients do not recognize levodopa-induced dyskinesias and do not experience any disability from the movements. As is, treating every dyskinesia is not necessarily essential, and clinicians focus on preventing worsening or reducing only disabling, bothersome dyskinesia with medical or surgical strategies (Aquino and Fox, 2015). Dyskinesia is most common at the peak-level of levodopa action, and consists of chorea, dystonia, and ballism, and to a lesser extent myoclonus (Nutt, 1990). Choreic movements in the limbs are the most common form of peak-dose dyskinesia, but dystonic posturing in the limbs can also occur. These involuntary movements may be initially mild and mainly involve the neck, and less commonly affect lips and jaw. They later spread to involve the trunk and can become more bothersome movements (Aquino and Fox, 2015). Myoclonus is a brief, involuntary, irregular twitching of a muscle or a group of muscles. Levodopa-induced myoclonus has been described as either spontaneous, action induced, or stimulus sensitive, and occurs within 10 to 20 minutes of levodopa administration (Aquino and Fox, 2015). Dyskinesia also occurs at low dose of levodopa action. Dyskinesia occurring during off-period is predominantly dystonic and mostly affects the legs and feet (Luquin et al., 1992). Off-period dystonia may completely disappear after withdrawal of levodopa for a few days or weeks (Aquino and Fox, 2015). When the levels of levodopa are rising and falling, at the beginning or end-of-dose respectively, dyskinesia can also occur and is known as “diphasic dyskinesia”. Diphasic dyskinesias are less common, tend to mainly affect the legs, and can involve slow stereotypical alternating leg movements (Luquin et al., 1992).

1.2 Impulse control disorders

Impulse control disorders include conditions involving problems in the self-control of emotions and behaviors such as pyromania or kleptomania (American Psychiatric Association, 2013). The fifth edition of the Diagnostic and Statistical Manual of Mental Disorders (American Psychiatric Association, 2013) has a specific chapter on disruptive, impulse-control and conduct disorders. Impulse control disorders have been studied among college students (Leppink et al., 2016b; Odlaug and Grant, 2010) and elderly patients (Tamam et al., 2014), as well as in several other disorders.

1.2.1 Definition of specific impulse control disorders

Intermittent explosive disorder The main feature of intermittent explosive disorder is recurrent behavioral outbursts representing a failure to control aggressive impulses. They are manifested by either verbal or physical aggression, or destruction or belongings. The magnitude of aggressiveness during these outbursts highly exceeds the provocation or any anticipated psychosocial stressors. The recurrent outbursts

are not purposeless, unpremeditated, and cause either impairment in occupational or interpersonal functioning, or marked distress in the individual ([American Psychiatric Association, 2013](#)).

Kleptomania Kleptomania is characterized by recurrent failures to resist impulse to steal objects that are not needed for their monetary value or for personal use. Individuals with kleptomania feel an increasing sense of tension shortly before committing the theft, then pleasure, relief or gratification at the time of commitment. Stealing is not committed in response to a hallucination or a delusion, or to express vengeance or anger. ([American Psychiatric Association, 2013](#))

Pyromania Pyromania is characterized by the multiple episodes of deliberate and purposeful fire setting. Persons with pyromania experience affective or tension arousal before setting a fire. They are interested in fire in situational contexts and are often regular watchers at fires. The fire setting is not done for profit and does not result from impaired judgement ([American Psychiatric Association, 2013](#)).

Pathological gambling The main feature of pathological gambling is recurrent and persistent dysfunctional patterns of gambling behavior leading to clinically significant distress or impairment. Pathological gamblers may need to gamble increasing amounts of money in a bid to achieve the desired excitement, may be irritable or restless when trying to cut down or stop gambling, often gamble when feeling distressed or lie to conceal the amount of involvement with gambling. The gambling behavior is not better explained by a manic episode. ([American Psychiatric Association, 2013](#)).

Compulsive sexual behavior Compulsive sexual behavior, also called hypersexuality, is characterized by persistently or recurrently present sexual or erotic thoughts or fantasies and desire for sexual activity ([American Psychiatric Association, 2013](#)). Individuals with this disorder feel driven or compelled to perform the behavior, which may cause distress ([Dell’Osso et al., 2006](#)). Impulsive-compulsive sexual disorder include unconventional sexual behaviors with a disturbance in the object of sexual gratification or in the expression of sexual gratification, and conventional sexual behaviors that have become excessive or uncontrolled ([Dell’Osso et al., 2006](#)).

Internet addiction disorder Internet addiction disorder is defined as a persistent and recurrent use of the Internet ([American Psychiatric Association, 2013](#)). Major symptoms of this disorder include preoccupation with the Internet, unsuccessful repeated efforts to decrease Internet use, staying online longer than intended, jeopardizing significant relationship, job, educational or career opportunity because of the Internet,

lying to relatives or physicians about the involvement with the Internet, and using the Internet a way of escaping from problems or of regulating mood (Young, 1998).

Compulsive buying disorder Compulsive buying disorder is characterized by excessive shopping cognitions and buying behavior that leads to distress or impairment (Black, 2007). Individuals with compulsive buying disorder are preoccupied with shopping and spending, and devote significant time to these behaviors. Shopping and spending are highly associated as window shopping is an uncommon pattern. Compulsive buying behaviors can be split into four phases: anticipation, preparation, shopping, and spending (Black, 2007). In the first phase, individuals develop preoccupations and thoughts with either shopping or having a specific item. The second phase consists in preparing for shopping and spending, such as deciding where to go. The third phase involving shopping and spending itself, which often procures high excitement and even sexual feelings (Schlosser et al., 1994). The fourth phase consists of the moment after the purchase, which is often experienced as a disappointment or a letdown. Common purchased items during these behaviors include clothing, shoes, and household items (Christenson et al., 1994; Miltenberger et al., 2003; Schlosser et al., 1994).

Binge eating disorder The main feature of binge eating disorder is recurrent episodes of binge eating. An episode of binge eating is characterized by eating, in a discrete span of time, a much larger amount of food than most people would eat under similar circumstances and in a similar period of time, and a sense of lack of control over eating during this span of time. These episodes are often associated with eating much more quickly than normal, eating until feeling uncomfortably full, eating large food quantities when not feeling physically hungry, eating alone because of embarrassment, and feeling disgusted with oneself, depressed, or very guilty afterward. (American Psychiatric Association, 2013)

Excoriation disorder Excoriation disorder, also called skin picking disorder, is characterized by recurrent skin picking resulting in skin lesions (American Psychiatric Association, 2013). Individuals with excoriation disorder experienced repeated attempts to curtail or stop skin picking. The skin picking causes significant distress or impairment in important areas of functioning and cannot be attributed to the physiological effects of any medical condition. Although it is now classified as an obsessive-compulsive disorder (American Psychiatric Association, 2013), excoriation disorder used to be classified in the impulse disorder category in the tenth revision of the International Statistical Classification of Diseases and Related Health Problems (Black and Grant, 2014).

1.2.2 Studies on impulse control disorders in subpopulations

Impulse control disorders have been studied in subpopulations, defined by an age range or another disorder.

Besides Parkinson's disease, which we will detail in the next section, the most studied disorders are obsessive-compulsive disorders (Fontenelle et al., 2005; Grant et al., 2006). ICDs and obsessive-compulsive disorders overlap in their phenomenology, co-morbidity, pathophysiology and family history (Fontenelle et al., 2011). Compulsive and impulsive disorders have been viewed at the opposite ends of a single dimension: the former motivated by a desire to avoid harm and the latter by reward-seeking behavior (Fineberg et al., 2010). The prevalence of ICDs in patients with obsessive-compulsive disorders is estimated to be around 11–35%, skin picking being the most common comorbid ICD (Fontenelle et al., 2005; Grant et al., 2006, 2010).

A review on ICDs and bipolar disorder highlighted the high number of common features for both disorders: phenomenological similarities, including pleasurable, dangerous, or harmful behaviors, impulsivity, and similar affective symptoms and dysregulation; onset in late childhood or early adulthood with episodic and/or chronic course; high comorbidity with one another and comparable comorbidity with other psychiatric disorders; high familial rates of mood disorder; and response to mood stabilizers and antidepressants (McElroy et al., 1996).

Impulse control disorders have been reported in Tourette syndrome (Jankovic and Kurlan, 2011), restless leg syndrome (Cornelius et al., 2010) and Perry syndrome (Mishima et al., 2015).

Impulse control disorders have also been studied in university students. In this subpopulation, the prevalence is estimated to be around 10% (Odlaug and Grant, 2010) and are associated with stress (Leppink et al., 2016b) and depression (Leppink et al., 2016a). Among elderly people, the prevalence is estimated to be around 17% and ICDs are associated with childhood conduct disorder and alcohol/substance abuse (Odlaug and Grant, 2010).

1.3 Impulse control disorders in Parkinson's disease

Since the first reports of impulse control disorders in Parkinson's disease in the early 2000s, impulse control disorders have been increasingly recognized. Given their potential impact on life functioning, including activities of daily living, interpersonal relationships, and social-occupational functioning, clinicians growingly pay specific attention to these impulsive behaviors (Weintraub and Claassen, 2017). ICDs are actually not symptoms of Parkinson's disease itself, but adverse effects of dopamine replacement therapy (de la Riva et al., 2014). ICDs have been broadly studied in PD, from prevalence to assessment to comorbidities.

1.3.1 Epidemiology

One of the earliest case reports dates back to 2003, which identified nine patients (0.5% of the sample) with pathological gambling (Driver-Dunkley et al., 2003). In 2010, the DOMINION study aimed at evaluating the point prevalence estimates of the four main ICDs among 3090 medicated PD patients in the United States and Canada (Weintraub et al., 2010a). One or more ICDs were identified in 13.6% of patients (compulsive buying in 5.7%, gambling in 5.0%, binge-eating disorder in 4.3%, and compulsive sexual behavior in 3.5%), with 3.9% of participants having two or more ICDs. A more recent longitudinal analysis of ICDs in a French research cohort estimated the 5-year cumulative incidence to be about 46% (Corvol et al., 2018).

As for any psychiatric disorder, environmental factors may influence the presence of ICDs in PD. In particular, cultural factors seem to impact the prevalence of specific ICDs. Studies in Turkey and India reported very low prevalences for pathological gambling (Kenangil et al., 2010; Sarathchandran et al., 2013), while this ICD has one of the highest prevalence in most Western studies (Baig et al., 2019; Garcia-Ruiz et al., 2014; Hurt et al., 2014; Weintraub et al., 2010a). Gambling is illegal in Turkey and heavily restricted in India, whereas it is legal in Western countries, and an important part of the American culture. Various studies lack uniformity to assess ICDs, and the definition itself of ICDs is subject to cultural differences (Weintraub and Claassen, 2017).

1.3.2 Assessment and diagnosis

As impulse control disorders have been increasingly recognized in Parkinson's disease, several screening tools and rating scales have been developed and used to assess and diagnose them.

The Movement Disorder Society-Sponsored Revision of the Unified Parkinson's Disease Rating Scale (MDS-UPDRS) has an item entitled *Features of dopamine dysregulation syndrome* in the part assessing non-motor aspects of experiences of daily living (Goetz et al., 2008). This single item encompasses impulse control disorders, dopamine dysregulation syndrome, punding, and hobbyism.

The Questionnaire for Impulsive-Compulsive Disorders in Parkinson's Disease (QUIP) was developed as a screening instrument for ICDs and related behaviors and is structured to be consistent with diagnostic criteria or defining clinical characteristics as described in the Diagnostic and Statistical Manual of Mental Disorders (Weintraub et al., 2009). Its three sections focus on (i) the four most common ICDs, (ii) punding and hobbyism, and (iii) compulsive medication use.

The Rating Scale version of the QUIP (QUIP-RS) was derived from the QUIP to measure the severity of ICDs (Weintraub et al., 2012). Each item is rated on a 5-point Likert scale and assesses the frequency of the symptoms with a range of scores from 0 (never) to 4 (very often). The sections are similar than the ones in the QUIP, although a

slight difference is that punding and hobbyism have been grouped together (Evans et al., 2019). This scale has been validated in several countries (Choi et al., 2020; Marques et al., 2019; Probst et al., 2014).

The Ardouin Scale of Behavior in Parkinson's Disease consists of eighteen items addressing non-motor symptoms, grouped in four parts: general psychological evaluation, apathy, non-motor fluctuations and hyperdopaminergic behaviors (Ardouin et al., 2009).

The Scale for Outcomes in Parkinson's Disease – Psychiatric Complications is a screening and severity scale that consists of a 7-item questionnaire (Visser et al., 2007). Two items are related to impulsive control disorders: one item for compulsive shopping and pathological gambling, and another one for hypersexuality. Each item score range from 0 (no symptoms) to 3 (severe symptoms) (Evans et al., 2019).

The Minnesota Impulsive Disorders Interview was originally developed in 2008 for the diagnosis of compulsive buying, trichotillomania, kleptomania, pyromania, intermittent explosive disorder, pathological gambling, and compulsive sexual behavior (Chamberlain and Grant, 2018; Grant, 2008). The original version was revised to match the changes made in the fifth edition of the Diagnostic and Statistical Manual of Mental Disorders (Chamberlain and Grant, 2018).

1.3.3 Associations

A wide range of factors have been associated with impulse control disorders in Parkinson's disease, from personality traits to psychiatric comorbidities to medications.

Demographics Significant differences between sexes have been observed, with men developing more pathological gambling and hypersexuality disorders and women developing more compulsive buying and eating disorders (Weintraub and Claassen, 2017). A younger age has been associated with ICDs in PD in numerous studies (Callesen et al., 2014; Poletti et al., 2013; Pontieri et al., 2015; Weintraub et al., 2010a). The DOMINION study, with 3090 PD patients from the United States and Canada, reported that PD patients with ICDs were most likely unmarried and living in the United States (Weintraub et al., 2010a).

Personality traits Unsurprisingly, the most assessed personality trait was impulsivity, with studies reporting higher impulsivity scores (Sáez-Francàs et al., 2016; Voon et al., 2011) and greater choice impulsivity (Sáez-Francàs et al., 2016). PD patients with ICDs were described as individuals with a higher level of neuroticism and lower levels of agreeableness and conscientiousness (Callesen et al., 2014), particularly among PD patients with pathological gambling (Gescheidt et al., 2016) or hypersexuality (Sachdeva et al., 2014). These patients were also reported to have ineffective coping skills (Olley et al., 2015).

Addictive disorders A family history of pathological gambling was significantly more prevalent among PD patients with ICDs than those without (Weintraub et al., 2010a). Past and current cigarette smoking has been associated with ICDs (Valença et al., 2013; Weintraub et al., 2010a), so was substance use of caffeine (Bastiaens et al., 2013; Gescheidt et al., 2016), tea, and alcohol (Ramírez Gómez et al., 2017).

Psychiatric comorbidities Mental illness was reported to be associated with the presence of ICDs, in particular anxiety and depression (Grall-Bronnec et al., 2018; Olley et al., 2015). Higher score of depression (Callesen et al., 2014; Joutsa et al., 2012; Vela et al., 2016; Voon et al., 2011), symptoms of depression (Gescheidt et al., 2016; Pontone et al., 2006), and a history of depression (Auyeung et al., 2011) have been found to be correlated with ICDs. Higher score of anxiety (Leroi et al., 2012; Pontieri et al., 2015; Sachdeva et al., 2014; Voon et al., 2011), trait anxiety (Sáez-Francàs et al., 2016), and a history of anxiety (Auyeung et al., 2011) were also found to be associated with ICDs. Only one study reported a higher obsessive-compulsive score (Voon et al., 2011), although both disorders share several features.

Sleep disturbances PD patients with ICDs have been reported to have an increased prevalence of sleep disturbances, including daytime sleepiness, worse sleep efficiency, and restless leg syndrome symptoms (Marques et al., 2018; O’Sullivan et al., 2011; Pontieri et al., 2015; Scullin et al., 2013). A strong association was shown between ICD symptoms, specially pathological gambling, and rapid eye movement sleep behavior disorder (Fantini et al., 2015, 2018, 2019, 2020; Ramírez Gómez et al., 2017).

Disease-related factors Several disease-related factors have been associated with ICDs. Younger age of PD onset (Callesen et al., 2014; Lee et al., 2010; Pontieri et al., 2015; Ye et al., 2011) and longer disease duration (Callesen et al., 2014; Lee et al., 2010; Pontieri et al., 2015) have been found correlated with ICDs. Association between ICDs and higher motor impairment has also been reported (Bastiaens et al., 2013; Callesen et al., 2014; Leroi et al., 2012). A negative association between motor fluctuations or dyskinesias and ICDs has been reported in one study (Ramírez Gómez et al., 2017). In particular, a higher score on the MDS-UPDRS Part I was found in two studies (Rodríguez-Violante et al., 2014; Sáez-Francàs et al., 2016), but it must be noted that one of the item of the MDS-UPDRS Part I is about dopamine dysregulation syndrome.

Medications Dopamine replacement therapy, specially dopamine agonists, has been strongly associated with ICDs. Ever use, longer cumulative duration, and higher cumulative dose of DAs have been correlated with ICDs (Corvol et al., 2018). The six dopamine agonists that have been approved by the US Food and Drug Administration (pramipexole, ropinirole, cabergoline, bromocriptine, rotigotine, and apomorphine) have

all been associated with ICDs (Grall-Bronnec et al., 2018). Dopamine agonists with a preferential affinity for D2-like receptors (D2 and D3 receptors), that is pramipexole and ropinirole, have been reported to have the strongest associations (Moore et al., 2014). To a lesser extent, associations with levodopa (Pontieri et al., 2015; Weintraub et al., 2010a) and amantadine (Weintraub et al., 2010b) have also been reported.

Genetic factors Association between ICDs and several single-nucleotide polymorphisms (SNPs) have been suggested in the following genes: *DRD3* (Castro-Martínez et al., 2018; Krishnamoorthy et al., 2016; Lee et al., 2009), *GRIN2B* (Lee et al., 2009; Zainal Abidin et al., 2015), *HTR2A* (Kraemmer et al., 2016; Lee et al., 2012), *ANKK1* (Hoenicka et al., 2015), *DRD1* (Erga et al., 2018; Zainal Abidin et al., 2015), *DRD2* (Kraemmer et al., 2016; Zainal Abidin et al., 2015), *OPRM1* (Cormier-Dequaire et al., 2018), *DAT1* (Cormier-Dequaire et al., 2018), *LRRK2* (Simuni et al., 2020), *GBA* (Simuni et al., 2020), *OPRK1* (Cormier-Dequaire et al., 2018; Kraemmer et al., 2016), and *SLC22A1* (Redenek et al., 2019). However a few studies did not report associations between ICDs and several SNPs from the following genes: *DRD2* (Cormier-Dequaire et al., 2018; Vallelunga et al., 2012), *COMT* (Vallelunga et al., 2012), *DAT1* (Vallelunga et al., 2012), *GRIN2B* (Cormier-Dequaire et al., 2018), and *HTR2A* (Cormier-Dequaire et al., 2018). The Parkinson's disease polygenic risk score has been reported not to be associated with ICDs (Ihle et al., 2020).

1.3.4 Prediction

While the literature on correlates with ICDs is large, studies focusing on the prediction of ICDs are very scarce: only three studies with a prediction task have been identified (Erga et al., 2018; Jesús et al., 2020; Kraemmer et al., 2016).

Kraemmer and others (Kraemmer et al., 2016) developed a clinical-genetic model to predict incident impulse control disorders in PD. The clinical features consisted of age, sex, PD treatment and duration of follow-up. The genetic variables consisted of thirteen candidate variants selected from the following genes: *DRD2-3*, *DAT1*, *COMT*, *DDC*, *GRIN2B*, *ADRA2C*, *SERT*, *TPH2*, *HTR2A*, *OPRK1*, and *OPRM1*. They worked on the Parkinson's Progression Markers Initiative (PPMI) database, which is an ongoing longitudinal multi-centre international study designed to identify biomarkers of PD progression in de novo and drug-naïve (at baseline) patients with PD. The algorithm trained with the aforementioned variables was a logistic regression.

Erga and others (Erga et al., 2018) also developed a clinical-genetic model, with slight differences compared to the previous study. The clinical features consisted only of age and PD treatment. The genetic variables consisted of fifty-six candidate variants selected from the following genes: *ADRA2C*, *DRD1-5*, *SLC6A3*, *DDC*, *COMT*, *SLC6A4*, *TPH2*, *HTR2A*, *OPRM1*, *OPRK1*, *GRIN2B*, and *BDNF*. The research cohort used was the Norwegian ParkWest study, which is a population-based longitudinal study of

incident PD. The trained algorithm was a logistic regression with an elastic-net penalty (Zou and Hastie, 2005).

Jesús and others (Jesús et al., 2020) also developed a clinical-genetic model. The clinical features consisted of sex, age, age at PD onset, years of disease evolution, DA equivalent daily dose, and levodopa equivalent daily dose. The genetic variables consisted of twenty genetic variants selected from the following genes: *DDC*, *DRD1*, *DRD2*, *DRD3*, *COMT*, *HTR2A*, *GRIN2B*, *TPH2*, *OPRM1*, *OPRK1*, *ADRA2C*, and *BDNF*. They worked on a research cohort from the Movement Disorder Clinic of the University Hospital Virgen del Rocío in Seville, Spain. They also trained logistic regression models with the aforementioned variables.

1.3.5 Other behavioral addictions

Although they are not impulse control disorders, several other related behaviors have been reported in Parkinson’s disease and can have a substantial negative impact on the quality of life of the patients. They are often referred to as *other behavioral addictions* or *related behaviors* in the literature and consist of dopamine dysregulation syndrome, punding and hobbyism (Weintraub and Claassen, 2017).

Dopamine dysregulation syndrome is characterized by the intake of large doses of dopaminergic drugs in excess of that required to control motor symptoms, endless requests to physicians for larger doses of dopamine replacement therapy or self-escalation of these medications without medical approval despite severe social destructive behaviors (Cilia et al., 2014). The prevalence is estimated to be around 34% in an advanced stage of PD (Cilia et al., 2014). A few cases have been reported in restless leg syndrome (Leu-Semenescu et al., 2009; Salas et al., 2009). A recent systematic review identified only nine case reports of dopamine dysregulation syndrome in non-Parkinson’s disease (Cartoon and Ramalingam, 2019).

Punding was first used to describe the behavior of people addicted to amphetamine (Rylander, 1972; Schiørring, 1981). Punding is a complex stereotyped behavior characterized by an intense fascination with repetitive manipulations of technical equipment, hoarding, grooming, continual handling, examining, and sorting common objects, pointless driving or walkabouts, and the engagement in extended monologues devoid of content (Evans et al., 2004). Punding behaviors often arise from particular habits or pastimes: people who regularly tinkered with technical objects are more likely to develop this kind of punding. Punding is also influenced by subject’s previous occupation: office workers and clerks may shuffle papers or fiddle purposelessly with computers while a seamstress may collect and arrange buttons (Spencer et al., 2011). A study reported the case of a 23-year-old Parkinsonian woman who developed unusual behaviors such as ceaseless sewing, disassembly and reassembly of phones, and coloring of drawings (El Otmani et al., 2015).

Hobbysim is defined as an excessive interest in one or several hobbies such as physical activity, artistic endeavor, do-it-yourself or gardening. For instance, a study reported a 77-year-old Parkinsonian man who started to show excessive hobbyism of painting four years after disease onset ([Matsuda et al., 2018](#)).

1.4 Machine learning

Machine learning is the process of automatically learning from data. Examples of tasks that machine learning can address include ([Hastie et al., 2009](#)):

- Predict the 10-year risk of future coronary heart disease.
- Estimate the amount of glucose in the blood of a diabetic person given the infrared absorption spectrum of that person's blood.
- Recognize the numbers in a handwritten ZIP code from a digitized image.
- Identify the risk factors for prostate cancer based on clinical and demographic variables.

In most scenarios, one has a target (outcome) measurement that one wants to predict from a set of features. The outcome can be quantitative (amount of glucose) or qualitative (presence or absence of a specific disease). Quantitative outcomes correspond to regression tasks, while qualitative outcomes correspond to classification tasks. One has a *training set* to train an algorithm and a *test set* (replication set) to evaluate its performance. These scenarios are *supervised learning* problems, because the learning process is supervised by the target.

The next sections introduce the notations used for the data and some of the most common machine learning algorithms.

1.4.1 Notations

Let n be the number of samples and m be the number of features. We consider data sets consisting of a $n \times m$ matrix \mathbf{X} representing the input data and a n vector representing the target data:

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix} = \begin{bmatrix} x_{11} & \dots & x_{1m} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nm} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

When longitudinal data is available, that is data at several time points, a time index is added:

$$\mathbf{X}^{(t)} = \begin{bmatrix} \mathbf{x}_1^{(t)} \\ \vdots \\ \mathbf{x}_n^{(t)} \end{bmatrix} = \begin{bmatrix} x_{11}^{(t)} & \dots & x_{1m}^{(t)} \\ \vdots & \ddots & \vdots \\ x_{n1}^{(t)} & \dots & x_{nm}^{(t)} \end{bmatrix}, \quad \mathbf{y}^{(t)} = \begin{bmatrix} y_1^{(t)} \\ \vdots \\ y_n^{(t)} \end{bmatrix}$$

and the input matrices and target vectors are available at several time points:

$$\left(\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(t)}, \dots, \mathbf{X}^{(T)} \right), \quad \left(\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(t)}, \dots, \mathbf{y}^{(T)} \right)$$

When only one sample is considered, the sample index is omitted. Likewise, when only one time point, the time index is omitted. The input data is thus the vector \mathbf{x} and the target data is y .

In the case of regression, y is a real number. In the case of classification, y is a single label. In particular, for binary classification, we consider that both classes are denoted $+1$ and -1 , that is $y \in \{-1, +1\}$.

The objective is to predict y given \mathbf{x} . The prediction is denoted \hat{y} . The most general formulation is:

$$\hat{y} = g(f(\mathbf{x}))$$

where f is the *decision function* and g is the final prediction. For regression tasks, g is the identity function, and the decision function is the final prediction.

1.4.2 Algorithms

Linear models

A linear model is a model that linearly combines the features:

$$f(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\beta} = \beta_0 + \sum_{j=1}^m \beta_j x_j$$

The vector $\boldsymbol{\beta}$ consists of:

- the intercept (constant) β_0 , and
- the coefficients $(\beta_1, \dots, \beta_m)$, where each coefficient β_j is associated to the feature x_j .

The vector $\boldsymbol{\beta}$ defines an hyperplane and $f(\mathbf{x})$ corresponds to the distance of \mathbf{x} to this hyperplane. A hyperplane is a subspace whose dimension is one less than that of the original space. For instance, in the two-dimensional case, a hyperplane is a line. In the the three-dimensional case, a hyperplane is a plane.

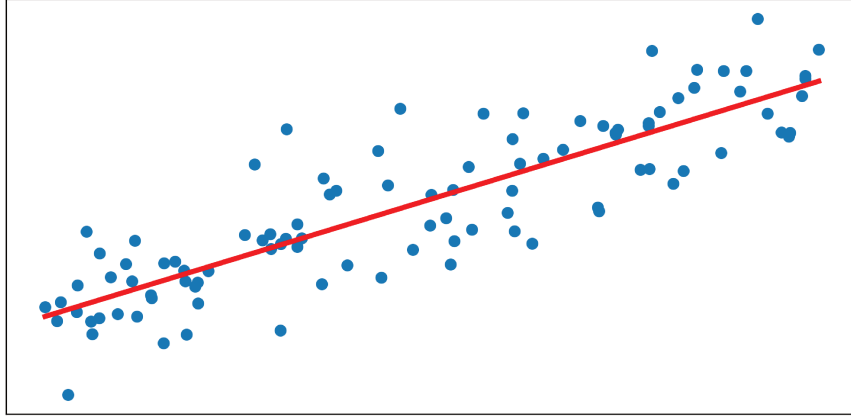


Figure 1.4: Ordinary least squares linear regression. When the input data consists of a single feature, the input and target variables can be visualized on a 2D plot. The ordinary least squares linear regression finds the best linear relationship (line) between both variables.

Ordinary Least Squares Linear Regression

The ordinary least squares linear regression is a linear regression model that is trained by minimizing the residual sum of squares between the observed targets in the data set, and the targets predicted by the linear approximation ([Hastie et al., 2009](#)):

$$\hat{y} = f(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\beta} = \beta_0 + \sum_{j=1}^m \beta_j x_j$$

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \min_{\boldsymbol{\beta}} \sum_{i=1}^n \left(y_i - \mathbf{x}_i^\top \boldsymbol{\beta} \right)^2 = \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$$

Figure 1.4 illustrates the main concept of the algorithm. When the input data consists of a single feature, the input and target variables can be visualized on a 2D plot. The ordinary least squares linear regression finds the best linear relationship between both variables. In this particular case, a linear relationship is simply a line.

Logistic Regression

For binary classification tasks, an hyperplane splits a space into two subspaces. In one subspace, the signed distance to the hyperplane is positive; in the other subspace, the signed distance to the hyperplane is negative. The decision of the algorithm depends on the sign of the signed distance:

$$\hat{y} = g(f(\mathbf{x})) = \begin{cases} +1 & \text{if } f(\mathbf{x}) > 0 \\ -1 & \text{if } f(\mathbf{x}) < 0 \end{cases}$$

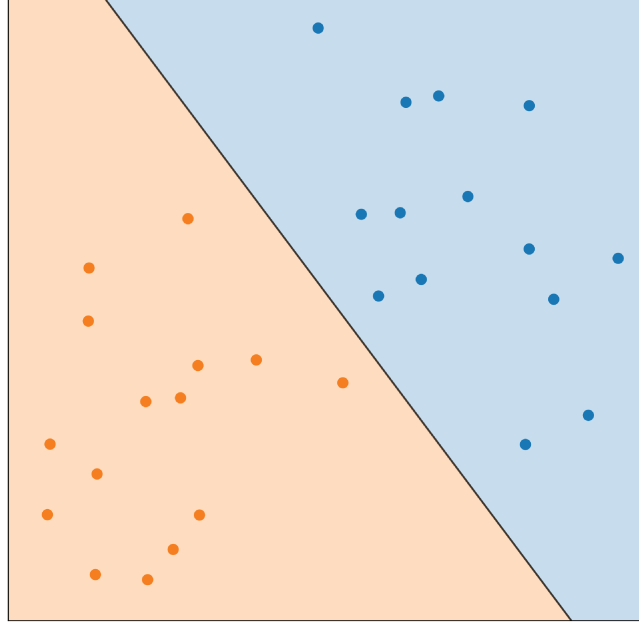


Figure 1.5: Decision function of a logistic regression model. A logistic regression is a linear model, that is its decision function is linear. In the two-dimensional case, it separates a plane with a line.

The logistic regression model transforms the signed distance to the hyperplane into a probability using the sigmoid function (Hastie et al., 2009):

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^m \beta_i x_i$$

$$\mathbb{P}(y = +1|\mathbf{x}) = \sigma(f(\mathbf{x})) = \frac{1}{1 + \exp(-f(\mathbf{x}))}$$

By applying the inverse of the sigmoid function, which is known as the logit function, one can see that the natural logarithm of the *odds ratio* is modeled as a linear combination of the features:

$$\log\left(\frac{\mathbb{P}(y = +1|\mathbf{x})}{\mathbb{P}(y = -1|\mathbf{x})}\right) = \log\left(\frac{\mathbb{P}(y = +1|\mathbf{x})}{1 - \mathbb{P}(y = +1|\mathbf{x})}\right) = f(\mathbf{x}) = \beta_0 + \sum_{j=1}^m \beta_j x_j$$

Figure 1.5 illustrates the decision function in the two-dimensional case where both classes are linearly separable.

Support Vector Machine

The original support vector machine (SVM) algorithm was invented in 1963 (Vapnik and Lerner, 1963). Figure 1.6 illustrates the main concept of this algorithm. When both classes are linearly separable, there exists an infinite number of hyperplanes that

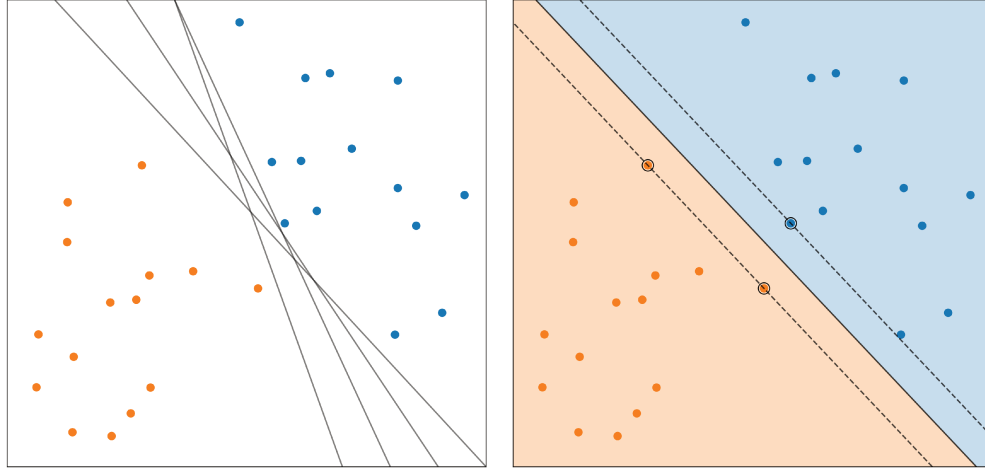


Figure 1.6: Support vector machine classifier. When two classes are linearly separable, there exists an infinite number of hyperplanes separating them (left). The decision function of the support vector machine classifier is the hyperplane that maximizes the margin, that is the distance between the hyperplane and the closest points to the hyperplane (right).

separate both classes. The SVM algorithm finds the hyperplane that maximises the distance between the hyperplane and the closest points of both classes to the hyperplane.

The SVM algorithm was extended in 1992 to non-linear decision functions using the *kernel trick* (Boser et al., 1992) and in 1995 to non-strictly separable classes (Cortes and Vapnik, 1995). The general form of this algorithm is:

$$f(\mathbf{x}) = \alpha_0 + \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$$

where k is the *kernel*. Popular kernels include:

- linear kernel: $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$
- polynomial kernel: $k(\mathbf{x}, \mathbf{x}') = (\gamma \mathbf{x}^\top \mathbf{x}' + c_0)^d$ with $c_0 \geq 0, d \in \mathbb{N}^*$
- sigmoid kernel: $k(\mathbf{x}, \mathbf{x}') = \tanh(\mathbf{x}^\top \mathbf{x}' + c_0)$ with $c_0 \geq 0$
- RBF kernel: $k(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|_2^2)$ with $\gamma > 0$

Figure 1.7 illustrates the decision functions for these kernels. Non-linear kernels allow for more complex decision functions. This is particularly useful when the data is not linearly separable, which is the most common use case.

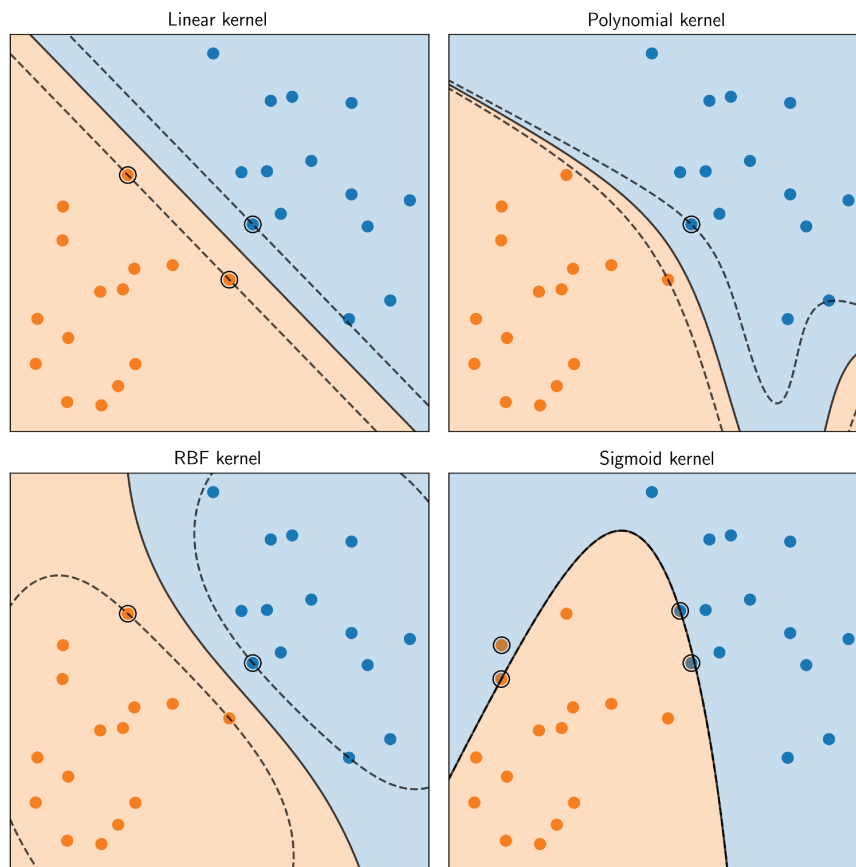


Figure 1.7: Impact of the kernel on the decision function of a support vector machine classifier. A non-linear kernel allows for a non-linear decision function.

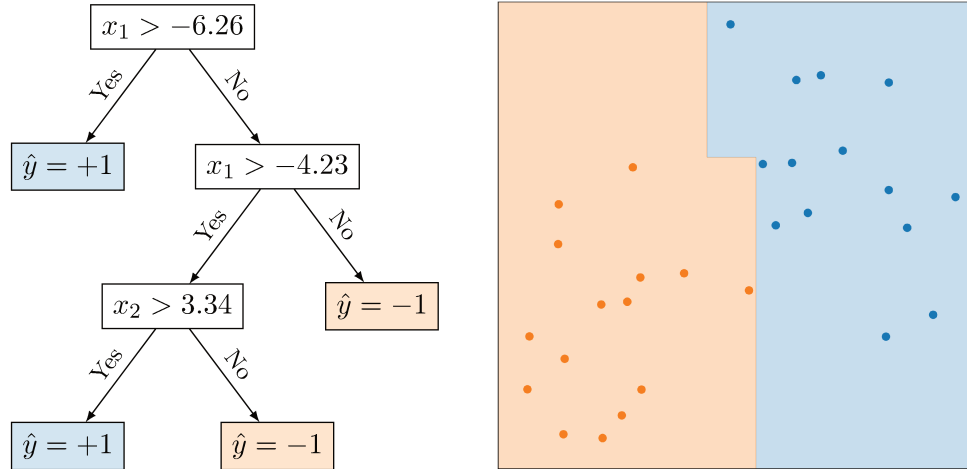


Figure 1.8: A decision tree: (left) the rules learned by the decision tree, and (right) the corresponding decision function.

Decision Tree

A decision tree is an algorithm containing only conditional statements and can be represented with a tree-like graph (Breiman et al., 1984). This graph consists of:

- decision nodes for the first condition,
- branches for the potential outcomes of each decision node, and
- leaf nodes for the final decision.

Figure 1.8 illustrates a decision tree and its corresponding decision function. For a given sample, the final decision is obtained by following its corresponding path, starting at the root node.

Random Forest

One limitation of decision trees is that they have a low bias but a high variance. An approach to overcome this limitation is to build an ensemble of trees. In order to have trees that are not perfectly correlated, subsets of the samples and the features are considered, introducing randomness. For each decision tree, only a subset of the samples are considered, usually drawn uniformly with replacement from the whole set. For each decision node of each tree, only a subset of the features are considered to find the best split. Both characteristics explain the name given to this algorithm: random forest (Breiman, 2001).

Gradient Tree Boosting

Each tree of a random forest is built independently from the other trees. This algorithm can easily take advantage of parallelization, which is an upside. An apparent limitation

is that each tree starts the training process all over again. Boosting is a technique that sequentially trains weak algorithms and sums them to obtain a strong algorithm (Breiman, 1996). Gradient boosting is a generalization of boosting by allowing optimization of an arbitrary loss function (Breiman, 1997). Gradient boosting algorithms can be seen as iterative functional gradient descent algorithms (Friedman, 2001; Mason et al., 2000). Gradient boosting is often used with decision trees, hence the name gradient tree boosting.

More specifically, a gradient boosting algorithm is the sum of weak algorithms:

$$f(\mathbf{x}) = f_0(\mathbf{x}) + f_1(\mathbf{x}) + \dots + f_H(\mathbf{x}) = \sum_{h=0}^H f_h(\mathbf{x})$$

and each weak algorithm is trained using functional gradient descent on the precedent weak algorithm:

$$f_h(\mathbf{x}) = f_{h-1}(\mathbf{x}) - \gamma_h \sum_{i=1}^n L(y_i, f_{h-1}(\mathbf{x}_i))$$

$$\gamma_h = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, f_{h-1}(\mathbf{x}_i) - \gamma \nabla_{f_{h-1}} L(y_i, f_{h-1}(\mathbf{x}_i)))$$

where L is the loss function that measures how good the predictions are. The lower, the better the predictions are, thus γ_h is the argument of the minimum.

Artificial Neural Network

Artificial neural networks are algorithms that can be represented by a network diagram as in Figure 1.9. They consist of a sequence of layers, allowing for extraction of high-level features from structured or unstructured data. A layer is often called an artificial neuron due to its similarity with a biological neuron. The artificial neuron receives one or more inputs and combines them to produce an output. The output is analogous to the axon of a biological neuron, and its value propagates to the input of the next layer, similarly to a synapse. Like electrical circuits, layers can be connected in series or in parallel, the former being much more common than the latter. The first layer is the input layer, consisting of the input data, and the final layer is the output layer, consisting of the prediction.

Several types of layers have been developed to deal with different types of data. Fully connected layers apply a linear transformation of the input followed by a non-linear activation function. Convolutions are commonly used for images and time series because each element of the input is strongly correlated to its neighbors. Recurrent units are dedicated to sequential data as they can take as input a variable number of elements. A typical application is natural language processing, because the number of

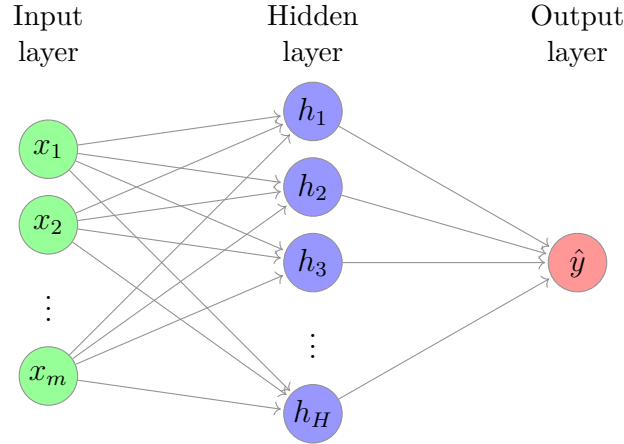


Figure 1.9: Example of an artificial neural network. This neural network has three layers: the input layer with the input features (x_1, \dots, x_m) , the hidden layer with the hidden features (h_1, \dots, h_H) extracted from the input features, and the output layer with the prediction \hat{y} made from the hidden features. This type of architecture is known as a multilayer perceptron.

words in a sentence is not constant. They can also be applied to longitudinal data.

Artificial neural networks typically have many parameters to be learned, with orders of magnitude ranging from thousands to billions of parameters. Training such algorithms is difficult as it requires a lot of data and computing power to estimate these parameters. Although research on artificial neural networks dates back to the late 1960s (Ivakhnenko and Lapa, 1967), their rise only occurred in the early 2010s. In 2012, the winning team of the ImageNet LSVRC-2010 contest used deep convolutional neural networks (Krizhevsky et al., 2012). The large increase of data and the constant progress in hardware broadened the applications of artificial neural networks and deep learning to many fields, including machine translation, object detection, image classification, chat bots, and so on.

As longitudinal cohorts naturally provide longitudinal data, the next section describes more precisely how recurrent neural networks work.

Recurrent neural networks

One of the key concepts of recurrent neural networks (RNNs) (Rumelhart et al., 1986) is sharing parameters across different parts of a model. Parameter sharing makes it possible to extend and apply the model to examples of different lengths, and generalize across them (Goodfellow et al., 2016).

A recurrent neural network is defined by the following recurrent equation

$$\mathbf{h}^{(t)} = f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \boldsymbol{\theta})$$

and is illustrated in Figure 1.10. The hidden state at time t , $\mathbf{h}^{(t)}$, is driven by:

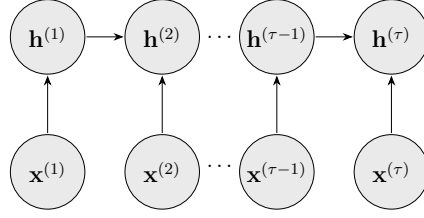


Figure 1.10: Main concept of a recurrent neural network. The hidden state at time t , $\mathbf{h}^{(t)}$, is derived from the hidden state at time $t - 1$, $\mathbf{h}^{(t-1)}$, and the external signal at time t , $\mathbf{x}^{(t)}$.

- the hidden state at time $t - 1$, $\mathbf{h}^{(t-1)}$,
- an external signal at time t , $\mathbf{x}^{(t)}$, and
- the parameters of the RNN, θ .

There is no time index for θ since θ is shared across the different time points. Intuitively, $\mathbf{h}^{(t)}$ represents the information extracted by the RNN at time t from $(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)})$ and is computed from:

- the information extracted by the RNN at time $t - 1$, $\mathbf{h}^{(t-1)}$, and
- the external signal at time t , $\mathbf{x}^{(t)}$,

which simply means that the extracted information is updated based on the new observation of the external signal.

The function f defines the relationship between $\mathbf{h}^{(t)}$, $\mathbf{h}^{(t-1)}$ and $\mathbf{x}^{(t)}$. Several functions have been developed and evaluated in the literature. We will briefly present the most used functions, often called *units*.

Vanilla Recurrent Neural Network The vanilla RNN unit was introduced by [Rumelhart et al. \(1986\)](#) and is illustrated in [Figure 1.11](#). The hidden state at time t , $\mathbf{h}^{(t)}$, is a linear combination of the hidden state at time $t - 1$, $\mathbf{h}^{(t-1)}$, and of the external signal at time t , $\mathbf{x}^{(t)}$, followed by an activation function, generally the hyperbolic tangent:

$$\mathbf{h}^{(t)} = \tanh \left(\mathbf{W}_{ih} \mathbf{x}^{(t)} + \mathbf{b}_{ih} + \mathbf{W}_{hh} \mathbf{h}^{(t-1)} + \mathbf{b}_{hh} \right)$$

One of the appeals of RNNs is their theoretical capability of connecting previous information to the present task. Unfortunately, the vanilla RNN unit is too simple to handle long dependencies in practice and suffers from the vanishing gradient problem.

Long Short-Term Memory The Long Short-Term Memory (LSTM) unit, introduced by [Hochreiter and Schmidhuber \(1997\)](#), was explicitly designed to avoid the long-term dependency problem and is illustrated in [Figure 1.12](#). The LSTM unit has a

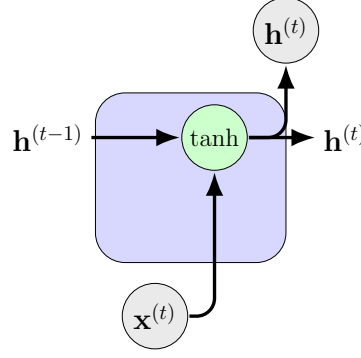


Figure 1.11: Vanilla recurrent neural network unit. The hidden state at time t , $\mathbf{h}^{(t)}$, is a linear combination of the hidden state at time $t - 1$, $\mathbf{h}^{(t-1)}$, and of the external signal at time t , $\mathbf{x}^{(t)}$, followed by an activation function.

cell state (\mathbf{C}) and four gates (\mathbf{f} , $\tilde{\mathbf{C}}$, \mathbf{i} and \mathbf{o}) that are updated at each time using the following equations:

$$\begin{aligned}
 \mathbf{f}^{(t)} &= \sigma \left(\mathbf{W}_{if} \mathbf{x}^{(t)} + \mathbf{b}_{if} + \mathbf{W}_{hf} \mathbf{h}^{(t-1)} + \mathbf{b}_{hf} \right) \\
 \mathbf{i}^{(t)} &= \sigma \left(\mathbf{W}_{ii} \mathbf{x}^{(t)} + \mathbf{b}_{ii} + \mathbf{W}_{hi} \mathbf{h}^{(t-1)} + \mathbf{b}_{hi} \right) \\
 \mathbf{o}^{(t)} &= \sigma \left(\mathbf{W}_{io} \mathbf{x}^{(t)} + \mathbf{b}_{io} + \mathbf{W}_{ho} \mathbf{h}^{(t-1)} + \mathbf{b}_{ho} \right) \\
 \tilde{\mathbf{C}}^{(t)} &= \tanh \left(\mathbf{W}_{ig} \mathbf{x}^{(t)} + \mathbf{b}_{ig} + \mathbf{W}_{hg} \mathbf{h}^{(t-1)} + \mathbf{b}_{hg} \right) \\
 \mathbf{C}^{(t)} &= \mathbf{f}^{(t)} \times \mathbf{C}^{(t-1)} + \mathbf{i}^{(t)} \times \tilde{\mathbf{C}}^{(t)} \\
 \mathbf{h}^{(t)} &= \mathbf{o}^{(t)} \times \tanh \left(\mathbf{C}^{(t)} \right)
 \end{aligned}$$

The cell state is the key component of the LSTM unit. The LSTM unit has the ability to add or remove information to the cell state, carefully regulated by the gates. The forget gate \mathbf{f} controls what the cell state must forget, while the input gate \mathbf{i} and the candidate gate $\tilde{\mathbf{C}}$ regulates the new information added to the cell state. Finally, the output gate \mathbf{o} controls which information of the cell state goes in the hidden state.

Gated Recurrent Unit The Gated Recurrent Unit was introduced by [Cho et al. \(2014\)](#) and is a variant of the LSTM unit with no output gate. The GRU has fewer parameters than the LSTM unit and has been reported to exhibit better performance on some smaller datasets ([Chung et al., 2014](#)). The GRU consists of three gates (\mathbf{r} , \mathbf{z}

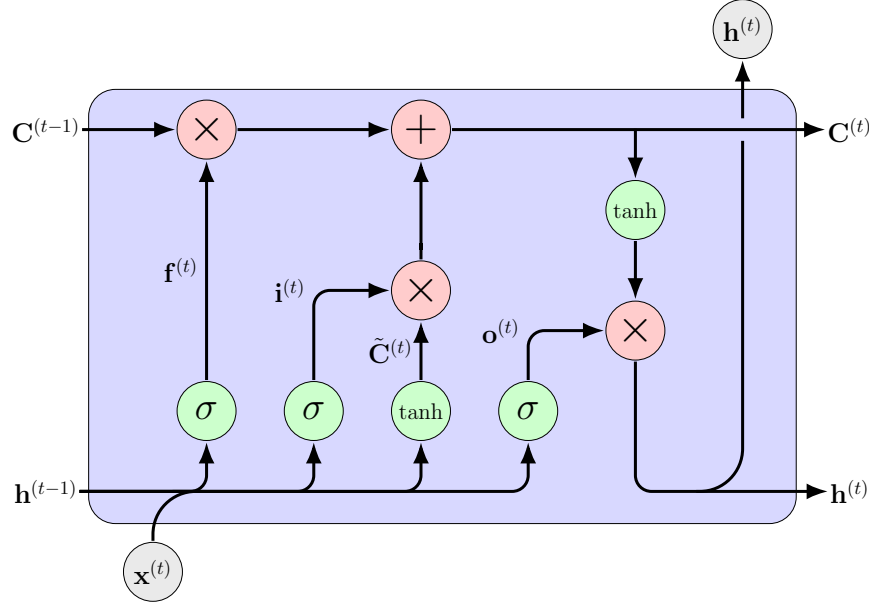


Figure 1.12: Long Short-Term Memory unit. The hidden state is updated via the the new observation, the cell state, and the four gates.

and $\tilde{\mathbf{h}}$) that are updated using the following equations:

$$\begin{aligned} \mathbf{z}^{(t)} &= \sigma \left(\mathbf{W}_{iz} \mathbf{x}^{(t)} + \mathbf{b}_{iz} + \mathbf{W}_{hz} \mathbf{h}^{(t-1)} + \mathbf{b}_{hz} \right) \\ \mathbf{r}^{(t)} &= \sigma \left(\mathbf{W}_{ir} \mathbf{x}^{(t)} + \mathbf{b}_{ir} + \mathbf{W}_{hr} \mathbf{h}^{(t-1)} + \mathbf{b}_{hr} \right) \\ \tilde{\mathbf{h}}^{(t)} &= \tanh \left(\mathbf{W}_{in} \mathbf{x}^{(t)} + \mathbf{b}_{in} + \mathbf{r}^{(t)} \times \left(\mathbf{W}_{hn} \mathbf{h}^{(t-1)} + \mathbf{b}_{hn} \right) \right) \\ \mathbf{h}^{(t)} &= \left(1 - \mathbf{z}^{(t)} \right) \times \mathbf{h}^{(t-1)} + \mathbf{z}^{(t)} \times \tilde{\mathbf{h}}^{(t)} \end{aligned}$$

The reset gate \mathbf{r} allows for resetting the hidden state with the observation \mathbf{x} , creating the candidate gate $\tilde{\mathbf{h}}$. The update gate \mathbf{z} allows for updating the hidden state with the candidate gate $\tilde{\mathbf{h}}$.

1.4.3 Regularization

Most machine learning algorithms are trained by minimizing a cost function:

$$\min_{\theta} c(\theta)$$

The cost function c measures the difference between the predictions of the algorithm and the true values of the target. The lower the loss function, the better the predictions. Fitting the training data consists in iteratively updating the parameters of the algorithm θ to minimize the loss function. However, if the model is too complex, its error on the training set is much lower than on the test, that is generalization (replication) on new

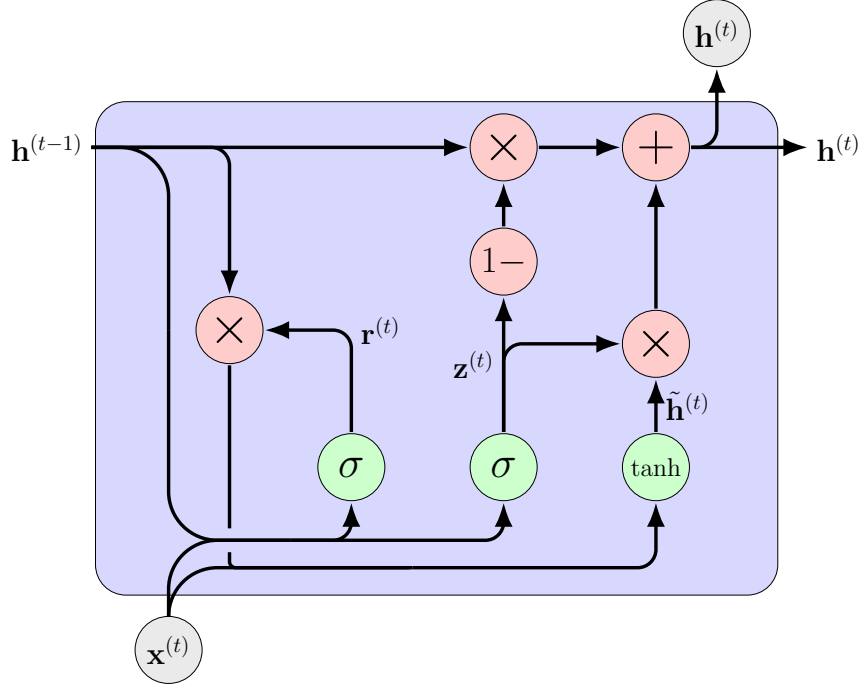


Figure 1.13: Gated Recurrent Unit. The Gated Recurrent Unit is a variant of the Long Short-Term Memory unit with no output gate (and therefore no cell state). The hidden state is updated via the new observation and the three gates.

observations is suboptimal. This phenomenon is known as *overfitting* (see Figure 1.14).

A common approach to avoid overfitting is to add a *regularization* term in the cost function that limits the complexity of the model:

$$\min_{\theta} c(\theta) + \lambda \times \text{Reg}(\theta)$$

The value of λ corresponds to the weight of the regularization in the loss function: the higher, the lower the complexity of the model. To illustrate the effect of regularization, we generate a toy data set from the following distributions (see Figure 1.15 for an example of a generated data set with 100 points):

$$\begin{aligned} x &\stackrel{\text{iid}}{\sim} \mathcal{U}_{[0,10]} \\ y &\sim \mathcal{N}(\sin(x), 0.5) \end{aligned}$$

A kernel ridge regression algorithm (Murphy, 2012) is trained on this data set for different values of λ (see Figure 1.16). When the value of λ is too high, the model does not fit the data enough (underfitting). When the value of λ is too low, the model fits the data too much (overfitting). An appropriate trade-off between fitting the data and limiting the complexity of the model gives the best results. This is known as the

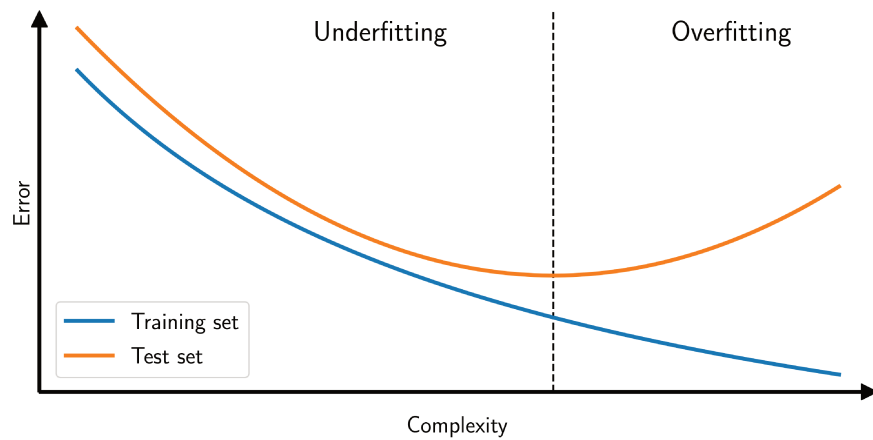


Figure 1.14: Relationship between error and model complexity. A too low complexity leads to underfitting. A too high complexity leads to overfitting.

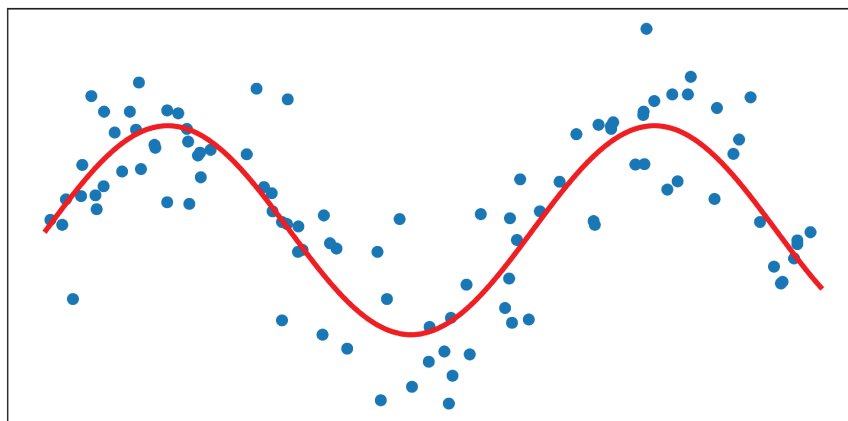


Figure 1.15: Toy regression data set with a non-linear relationship.

bias-variance trade-off. If the loss function is the squared difference between the true value y (fixed) and the predicted value \hat{y} (random variable), then its expected value is the sum of the squared bias of \hat{y} and its variance:

$$\begin{aligned}
\mathbb{E}[(y - \hat{y})^2] &= \mathbb{E}[y^2 - 2y\hat{y} + \hat{y}^2] \\
&= y^2 - 2y\mathbb{E}[\hat{y}] + \mathbb{E}[\hat{y}^2] \\
&= y^2 - 2y\mathbb{E}[\hat{y}] + \mathbb{E}[\hat{y}^2] + \mathbb{E}[\hat{y}]^2 - \mathbb{E}[\hat{y}]^2 \\
&= (\mathbb{E}[\hat{y}] - y)^2 + \mathbb{E}[\hat{y}^2] - \mathbb{E}[\hat{y}]^2 \\
&= (\mathbb{E}[\hat{y}] - y)^2 + \mathbb{E}[\hat{y}^2 - \mathbb{E}[\hat{y}]^2] \\
&= (\mathbb{E}[\hat{y}] - y)^2 + \mathbb{E}[\hat{y}^2 - 2\mathbb{E}[\hat{y}]^2 + \mathbb{E}[\hat{y}]^2] \\
&= (\mathbb{E}[\hat{y}] - y)^2 + \mathbb{E}[\hat{y}^2 - 2\hat{y}\mathbb{E}[\hat{y}] + \mathbb{E}[\hat{y}]^2] \\
&= (\mathbb{E}[\hat{y}] - y)^2 + \mathbb{E}[(\hat{y} - \mathbb{E}[\hat{y}])^2] \\
\mathbb{E}[(y - \hat{y})^2] &= \underbrace{(\mathbb{E}[\hat{y}] - y)^2}_{\text{bias}^2} + \underbrace{\mathbb{V}[\hat{y}]}_{\text{variance}}
\end{aligned}$$

When the model does not capture the regularities of the data, its bias is high but its variance is low (underfitting). When the model captures the noise of the data, its bias is low but its variance is high (overfitting).

The most common regularization terms for structured (tabular) data are the ℓ_2 -penalty, ℓ_1 -penalty, and the elastic net.

ℓ_2 -penalty

The ℓ_2 -penalty of $\boldsymbol{\theta}$ is the squared ℓ_2 -norm of $\boldsymbol{\theta}$, that is the sum of the squared elements in $\boldsymbol{\theta}$:

$$\ell_2(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_2^2 = \sum_i \theta_i^2$$

Adding this term to the loss function has several advantages: (i) it makes the optimization problem strictly convex; (ii) it forces each value θ_i not to be too large; (iii) for linear models, the coefficients become more robust to collinearity. Linear regression with ℓ_2 -penalty is commonly known as ridge regression ([Tikhonov et al., 1977](#)):

$$\min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2$$

ℓ_1 -penalty

The ℓ_2 -penalty forces the values of the parameters not to be too large, but does not incentive to make small values tend to 0. Indeed, the square of a small value is even smaller. When the number of features is large, or when interpretability is important,

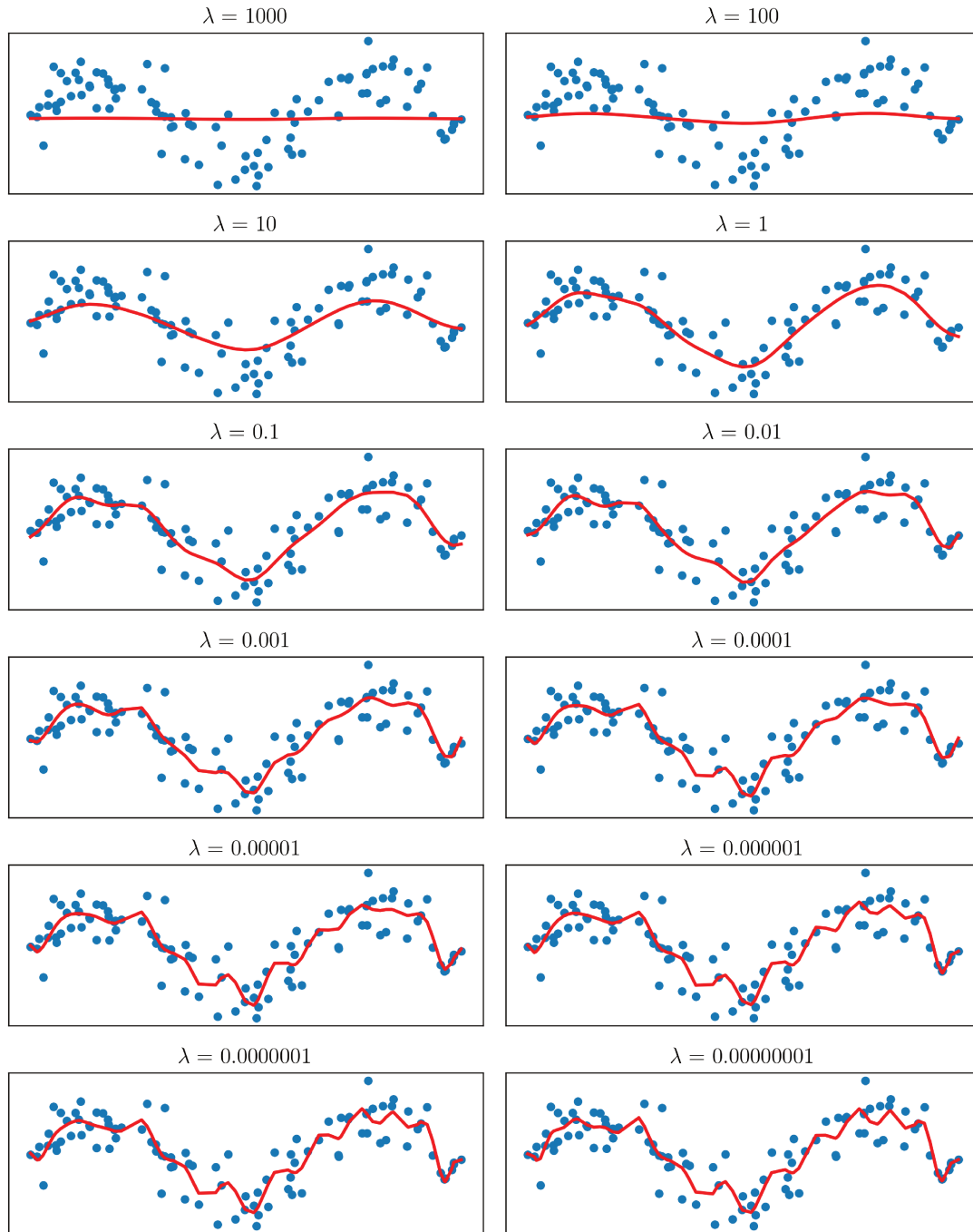


Figure 1.16: Illustration of regularization. A kernel ridge regression algorithm is fitted on this data with different values of λ , which is the weight of the regularization in the loss function. The smaller values of λ , the smaller the weight of the ℓ_2 regularization. The algorithm underfits (respectively overfits) the data when the value of λ is too large (respectively low).

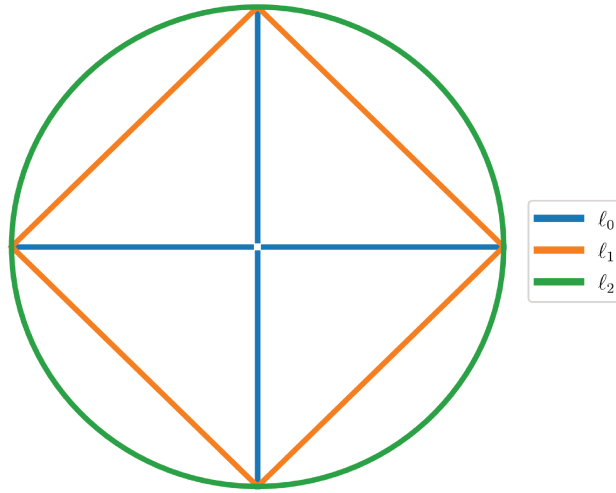


Figure 1.17: Unit balls of the ℓ_0 , ℓ_1 and ℓ_2 norms. For each norms, the set of points in \mathbb{R}^2 whose norm is equal to 1 is plotted. The ℓ_1 norm is the best convex approximation to the ℓ_0 norm. Note that the lines for the ℓ_0 extend to $-\infty$ and $+\infty$, but are cut for plotting reasons.

it can be useful to make the algorithm select the most important features. The corresponding norm is the ℓ_0 -norm, defined as the number of nonzero elements:

$$\ell_0(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_0 = \sum_i \mathbf{1}(\theta_i \neq 0)$$

However, the ℓ_0 -norm is not differentiable and not convex. The best convex approximation of the ℓ_0 -norm is the ℓ_1 -norm (see [Figure 1.17](#)), defined as the sum of the absolute values of each element:

$$\ell_1(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_1 = \sum_i |\theta_i|$$

Linear regression with ℓ_1 -penalty is commonly known as LASSO ([Tibshirani, 1996](#)):

$$\min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1$$

Elastic net

The ℓ_1 -penalty has several limitations. When the number of features is larger than the number of samples, the LASSO selects at most n variables before it saturates. When there is a group of highly correlated variables, the LASSO tends to select only one variable from a group and ignore the others ([Zou and Hastie, 2005](#)).

To overcome these limitations, the elastic net penalty linearly combines the ℓ_2 - and

ℓ_1 -penalties to get the best of both penalties (Zou and Hastie, 2005):

$$\text{EN}(\boldsymbol{\theta}) = \lambda \|\boldsymbol{\theta}\|_2^2 + \nu \|\boldsymbol{\theta}\|_1$$

1.4.4 Metrics

Metrics are a key component of machine learning as they evaluate how well models perform. There exist many ways of evaluating models, and choosing metrics must be based on the most important evaluation criteria, which depend on the task. In this section we will present the two metrics that we chose to evaluate the algorithms: the area under the receiver operating characteristic curve and the average precision. But first, we provide some reminder on binary classification and introduce the confusion matrix.

Confusion matrix

For classification tasks, a confusion matrix is matrix that reports all the possible combinations between the predicted output and the true output. Each row consists of the true classes, while each column consists of the predicted classes. For binary classification tasks, the four entries are known as:

- True positives (TP): the true class and the predicted class are both positive;
- False positives (FP): the true class is negative, and the predicted class is positive;
- True negatives (TN): the true class and the predicted class are both negative;
- False negatives (FN): the true class is positive, and the predicted class is negative.

Table 1.1 illustrates the concept of the confusion matrix for binary classification tasks, and the most common statistics derived therefrom. Since most of these statistics have several names, we recall them to avoid any confusion:

$$\begin{aligned} \text{True positive rate} = \text{Sensitivity} = \text{Recall} &= \frac{TP}{TP + FN} \\ \text{True negative rate} = \text{Specificity} &= \frac{TN}{TN + FP} \\ \text{Positive predictive value} = \text{Precision} &= \frac{TP}{TP + FP} \\ \text{Negative predictive value} &= \frac{TN}{TN + FN} \end{aligned}$$

The *true positive rate* (TPR), also known as *sensitivity* or *recall*, is the proportion of true positives among all the positives. The *true negative rate* (TNR), also known as *specificity*, is the proportion of true negatives among all the negatives. The *positive predictive value* (PPV), also known as *precision*, is the proportion of true positives

		Predicted classes		
		Positive	Negative	
True classes	Positive	TP	FN	$TPR = \frac{TP}{TP + FN}$
	Negative	FP	TN	$TNR = \frac{TN}{TN + FP}$
		$PPV = \frac{TP}{TP + FP}$	$NPV = \frac{TN}{TN + FN}$	

Table 1.1: Confusion matrix for binary classification. Each row represents the true positive and negative classes. Each column represents the predicted positive and negative classes. The entries of the confusion matrix correspond to all the possible outcomes. TP : True positives; FN : False negatives; FP : False positives; TN : True negatives; TPR : True positive rate; TNR : True negative rate; PPV : Positive predictive value; NPV : Negative predictive value.

among all the predicted positives. The *negative predictive value* (NPV) is the proportion of true negatives among all the predicted negatives.

A perfect classifier is classifier with no error, that is such that $FN = FP = 0$. For instance, a classifier is perfect if and only if:

- $TPR = TNR = 1$ since $TPR = 1 \iff FN = 0$ and $TNR = 1 \iff FP = 0$
- $TPR = PPV = 1$ since $TPR = 1 \iff FN = 0$ and $PPV = 1 \iff FP = 0$

Except if the binary classification task is relatively easy, having a perfect classifier is extremely rare. Sometimes, the *scores* of a classifier, such as probabilities, are as important as the predicted classes. For instance, clinicians are often more interested in the risks (i.e. probabilities) of a given disorder, rather than just a prediction. As a reminder, most binary classification algorithms consist of two steps:

1. Computing a score $f(\mathbf{x})$ for sample \mathbf{x}
2. Deriving the predicted class from the score $f(\mathbf{x})$ using a threshold ϵ :

$$\hat{y} = \begin{cases} +1 & \text{if } f(\mathbf{x}) > \epsilon \\ -1 & \text{if } f(\mathbf{x}) < \epsilon \end{cases}$$

For linear classifiers, $f(\mathbf{x}) = \beta_0 + \sum_{j=1}^m \beta_j x_j$ represents the signed distance to the hyperplane. Instead of comparing the true classes and the predicted classes

$$\mathbf{y} = (y_1, \dots, y_n) \quad , \quad \hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_n)$$

one can compare the true classes and the scores

$$\mathbf{y} = (y_1, \dots, y_n) \quad , \quad f(\mathbf{X}) = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$$

Two popular metrics to compare the true classes and the scores are the *area under the receiver operating characteristic curve* (ROC AUC), derived from the receiver operating characteristic (ROC) curve, and the *average precision* (AP) score, derived from the precision-recall (PR) curve.

Receiver operating characteristic curve

The receiver operating characteristic curve is the plot representing the true positive rate against the false positive rate at various threshold settings. The false positive rate is the proportion of false positives among all the negatives and can be calculated as $(1 - \text{specificity})$. The ROC curve starts at $(0, 0)$, when all the samples are predicted negatives ($\text{TPR} = 1 - \text{TNR} = 0$), and ends at $(1, 1)$, when all the samples are predicted positives ($\text{TPR} = 1 - \text{TNR} = 1$). The ROC curve is non-decreasing because no sample which is classified as a negative will ever be classified as a positive for any higher threshold.

The area under the ROC curve summarizes the ROC curve with a single score. ROC AUC can be computed as

$$\text{ROC AUC} = \frac{1}{n_{+1} \times n_{-1}} \sum_{\substack{i=1 \\ y_i=+1}}^n \sum_{\substack{j=1 \\ y_j=-1}}^n \mathbb{1}(f(\mathbf{x}_i) > f(\mathbf{x}_j))$$

where n_{+1} and n_{-1} are the number of positives and negatives respectively. ROC AUC has a simple interpretation: it is the probability that the classifier ranks a randomly chosen positive sample higher than a randomly chosen negative one. ROC AUC has the following properties:

- It is always between 0 and 1;
- The higher, the better;
- $\text{ROC AUC} = 0$ if and only if all the negative samples have higher scores than all the positive samples;
- $\text{ROC AUC} = 1$ if and only if all the positive samples have higher scores than all the negative samples (i.e. there exists a threshold yielding a perfect classifier);
- The expected ROC AUC of random guess is 0.5, independently of the distribution of the classes.

Precision-recall curve

The precision-recall curve is the the plot representing the precision against the recall at various threshold settings. Since the recall is the true positive rate, the only difference with the ROC curve is the replacement of the false positive rate with the precision. The PR curve starts at $(0, 1)$, when all the samples are predicted negatives ($\text{TPR} = 0, \text{PPV} = 1$), and ends at $(1, p)$, when all the samples are predicted positives ($\text{TPR} = 1, \text{PPV} = p$), p being the prevalence of the positive class. Precision is actually ill-defined when all the samples are predicted negatives (both the numerator and denominator are equal to 0), but the precision is expected to tend to 1 when the number of predicted positive samples tend to 0.

Contrary to the ROC curve, which takes into account the four possible outcomes (TP, FN, TN, FP), the precision-recall curve does not take into account the true negatives. Intuitively, precision is the ability of the classifier not to label as positive a sample that is negative, and recall is the ability of the classifier to find all the positive samples. The PR curve measures how well the classifier finds all the positive samples without labeling negative samples as positives.

The precision-recall curve is particularly useful for rare event detection, where false negatives are much more serious than false positives. For instance, let's consider a low-cost, fast, non-invasive diagnostic test for a serious disease. If this test is positive, the subject have a more advanced, costly, invasive exam to confirm the diagnosis. If this test is negative, the subject will have another test in a few years. We would like this test to find all the subjects with this disease while labeling the least healthy subjects as ill. Finding the healthy subjects is of much smaller interest than finding the ill ones.

The average precision score summarizes a precision-recall curve as the weighted mean of precisions achieved at each threshold, with the increase in recall from the previous threshold used as the weight:

$$\text{AP} = \sum_k (R_k - R_{k-1}) P_k$$

where P_k and R_k are respectively the precision and recall at the k -th threshold. Average precision is also the area under the precision-recall curve, but computed with a different technique than the one used to compute ROC AUC. AP is the area under the PR curve computed using the Riemann integral, while ROC AUC is the area under the ROC curve computed using the trapezoidal rule. A linear interpolation (with the trapezoidal rule) of points on the precision-recall curve provides an overly-optimistic measure of classifier performance (Davis and Goadrich, 2006; Flach and Kull, 2015). AP has the following properties:

- It is always between 0 and 1;
- The higher, the better;

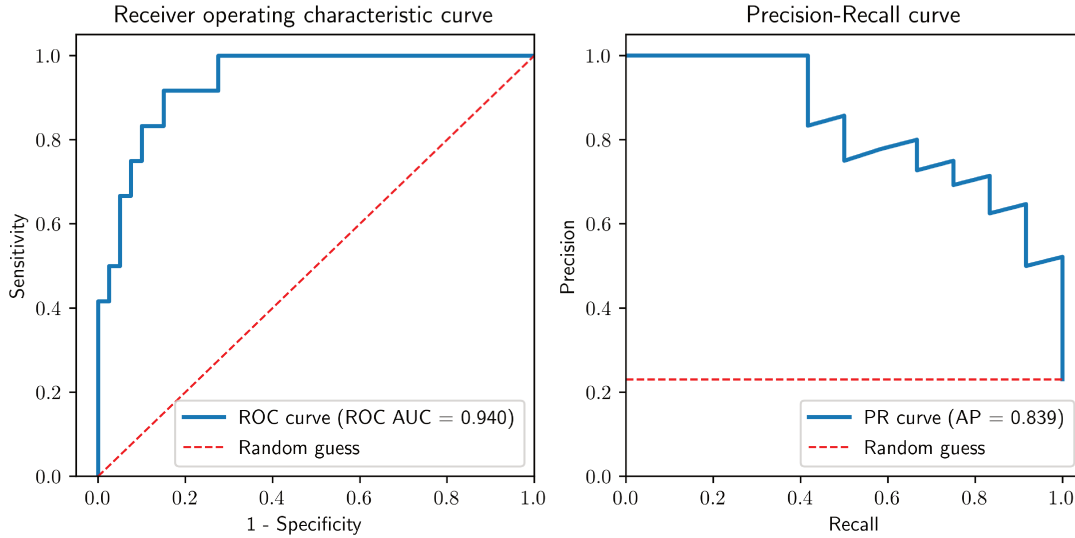


Figure 1.18: Receiver operating characteristic and precision-recall curves. The curves are represented in blue, and the expected curves of random guess are highlighted in red. The area under the ROC curve is computed with the trapezoidal rule, while the average precision score is the area under the PR curve computed with the Riemann integral.

- $AP = 1$ if and only if all the positive samples have higher scores than all the negative samples (i.e. there exists a threshold yielding a perfect classifier);
- The expected AP of random guess is equal to p , the prevalence of the positive class.

Figure 1.18 illustrates the ROC and PR curves. On the ROC curve, the lower the threshold, the higher the true positive and negative rates. The PR curve does not have a monotonicity property, as precision can increase when recall increases. Expected curves of random guess are highlighted.

1.5 Putting it all together

Given the association between impulse control disorders in Parkinson's disease and potential medical, financial, and/or legal medical complications, worse quality of life and function, high rates of psychiatric comorbidity, and strained interpersonal relationships and caregiver burden, identification and treatment of the symptoms are typically imperative (Weintraub and Claassen, 2017). Accurate prediction of ICDs ahead of time may allow avoiding their potentially terrible consequences and would enable to keep prescribing dopamine agonists with less fear of their adversarial effects.

Most of the literature on ICDs in PD is focused on *associations*, that is features that are correlated with ICDs, in cross-sectional studies. These associations, although important, have important limitations:

- Since these associations are cross-sectional, they may not hold true to predict ICDs *in advance*;
- Associations may be correlated, so that finding a new association may not be helpful to predict ICDs;
- Even new uncorrelated associations may not be helpful to predict ICDs.

Knowing associated factors may not be enough to prevent the irreversible consequences of ICDs. One of the earliest studies on pathological gambling in PD reported losses of hundreds of thousands of dollars for several patients (Gschwandtner et al., 2001). Being able to predict ICDs in advance could help preventing these life-changing events. A couple of studies focused on predicting ICDs, but the lack of cross-validation, that is evaluating a model on different data than the data used to train the algorithm, affects the confidence in their results.

Being able to accurately predict ICDs in advance would be of great interest clinically. Managing ICDs is not an easy task, as patients may not be aware of their behaviors, or can lie to their relatives and clinicians. Raising awareness in PD patients and their relatives would be more impactful if they could be told precisely when ICDs will be developed. ICDs may even be prevented, as case reports suggest that ICDs often resolve after reducing the dose of the existing dopamine agonists, in particular with complete discontinuation of DA treatment (Weintraub and Claassen, 2017).

Moreover, the genetic factors of ICDs in PD are mostly unknown. A few associations from candidate gene analyses have been reported, but these studies have not been replicated. Using machine learning algorithms to learn genetic factors from the whole genome, instead of picking a few genes based on prior information, may shed some new light on these genetic factors.

The interaction between putative genetic factors and clinical risk factors is also unknown, and may play an important role for the onset of ICDs in PD. Investigating different approaches modelling different types of interaction may help better understanding these disorders.

The association between known factors and impulse control disorders is complex (Grall-Bronnec et al., 2018). Impulse control disorders have been mainly studied from the statistical point of view, and machine learning has been underused. Machine learning, by automatically extracting information from data, could improve the predictability of ICDs and leverage the knowledge on this topic, but also improve the quality of life of patients and decrease caregiver burden.

1.6 Materials

Data analysis requires two essential components: data sets, to provide data, and software, to analyze data. Data sets can be private or publicly available, the former being

much more common than the latter in medicine. Software, by providing implementations of algorithms and utility tools, are at the core of data analysis.

In this section, we present the data sets from which we obtained data, and the software that we used to perform the analysis of these data.

1.6.1 Data sets

Medical data sets are not easy to collect and share for ethical, legal and privacy reasons. Nonetheless, data sets are needed to better understand disorders and to further advance scientific knowledge. In particular, publicly available data sets are very relevant because they give anyone access to data, regardless of their institution or employer. They also allow for tackling the reproducibility crisis that science is currently facing. Their main drawback is their possible overuse, leading to biases in the results.

In our work, we used two data sets: the Parkinson’s Progressive Markers Initiative (PPMI) database and the Drug Interaction With Genes in Parkinson’s Disease (DIGPD) database.

Parkinson’s Progressive Markers Initiative

In the field of Parkinson’s disease therapeutics, the ultimate goal is to develop disease-modifying treatments that slow, prevent or even reverse the underlying disease process. Validated biomarkers of disease progression would dramatically accelerate PD therapeutics research. However, current progression biomarkers are not optimal and are not fully validated.

The Parkinson’s Progression Markers Initiative (www.ppmi-info.org) is a landmark observational clinical study to comprehensively evaluate cohorts of significant interest using advanced imaging, biologic sampling and clinical and behavioral assessments to identify biomarkers of Parkinson’s disease progression. PPMI is taking place at clinical sites in the United States, Europe, Israel, and Australia (see [Table 1.2](#) for the full list of clinical sites).

Data and samples acquired from study participants enable the development of a comprehensive Parkinson’s database and biorepository, which is currently available to the scientific community to conduct field-changing research. PPMI follows standardized data acquisition protocols to ensure that tests and assessments conducted at multiple sites and across multiple cohorts can be pooled in centralized databases and repositories. The clinical, imaging and biologic data is easily accessible to researchers in real time through their website.

Drug Interaction With Genes in Parkinson’s Disease

The Drug Interaction With Genes in Parkinson’s Disease study is a longitudinal cohort study of patients with PD consecutively recruited from May 2009 to July 2013 in 4

Location	Organization
Athens, GREECE	National and Kapodistrian University of Athens
Atlanta, GA	Emory University
Baltimore, MD	Johns Hopkins University
Barcelona, SPAIN	Hospital Clinical de Barcelona
Birmingham, AL	University of Alabama at Birmingham
Boca Raton, FL	PD and Movement Disorders Center of Boca Raton
Boston, MA	Boston University
Chicago, IL	Northwestern University
Cincinnati, OH	University of Cincinnati
Cleveland, OH	Cleveland Clinic Foundation
Houston, TX	Baylor College of Medicine
Innsbruck, AUSTRIA	Innsbruck University
Kassel, GERMANY	Paracelsus-Elena Clinical Kassel / University of Marburg
London, UK	Imperial College London
New Haven, CT	Institute for Neurodegenerative Disorders
New York, NY	Columbia University Medical Center
New York, NY	Beth Israel Medical Center
Paris, France	Pitié-Salpêtrière Center
Philadelphia, PA	University of Pennsylvania
Portland, OR	Oregon Health & Science University
Rochester, NY	University of Rochester
Salerno, ITALY	University of Salerno
San Diego, CA	University of California, San Diego
San Francisco, CA	University of California, San Francisco
San Sebastian, SPAIN	Hospital Universitario Donostia
Seattle, WA	University of Washington
Sun City, AZ	Arizona Parkinson's Disease Consortium
Sunnyvale, CA	The Parkinson's Institute & Clinical Center
Sydney, AUSTRALIA	Macquarie University
Tampa, FL	University of South Florida
Tel Aviv, ISRAEL	Tel Aviv Sourasky Medical Center
Trondheim, NORWAY	Norwegian University of Science and Technology
Tübingen, GERMANY	Universität Tübingen

Table 1.2: PPMI clinical sites.

GA: Georgia; MD: Maryland; AL: Alabama; FL: Florida; IL: Illinois; OH: Ohio; TX: Texas; UK: United Kingdom; CT: Connecticut; NY: New York; PA: Pennsylvania; OR: Oregon; CA: California; WA: Washington; AZ: Arizona.

French university hospitals and 4 general hospitals (Corvol et al., 2018). Eligible patients were patients with PD (UK Parkinson’s Disease Society Brain Bank criteria) with disease duration shorter than 5 years at recruitment. After the baseline visit, annual clinical evaluations were performed over 5 years by movement disorders specialists who checked whether patients still fulfilled UK Parkinson’s Disease Society Brain Bank criteria at each visit and filled out standardized questionnaires. All patients had a blood sampling for DNA extraction and genome-wide genotyping. The study was conducted according to Good Clinical Practice Guidelines, and sponsored by Assistance Publique Hôpitaux de Paris. All patients provided informed consent, and the study was approved by local ethical committee and regulation authorities.

Assessed phenotypes

Parkinson’s disease is characterized by a wide range of symptoms. In order to diagnose them and assess their severity, screening tools and rating scales are used. General screening tools are often administered, but questionnaires and scales specific to Parkinson’s disease have also been developed.

For most phenotypes, several screening tools are available. Apart from the Movement Disorder Society-sponsored revision of the Unified Parkinson’s Disease Rating Scale, which quantifies the severity of Parkinson’s disease, there is no real consensus on the scales to use, and different databases may use different scales to assess some phenotypes. These differences can cause some issues for machine learning algorithms that expect the same features in the sets used to train and evaluate the models. Table 1.3 lists the questionnaires and scales used in PPMI and DIGPD for the main symptoms of Parkinson’s disease. Several symptoms that have been associated with impulse control disorders, such as anxiety and depression, are assessed with different scales in both databases.

1.6.2 Software

We describe in this section the software that we used, grouped by programming languages: Python for machine learning, C/C++ for genetic analyses, and R for meta-analysis.

Python

Python (Van Rossum and Drake, 2009) is a programming language that is easy to pick up regardless of past programming experience. Python is developed under an open source license, making it freely usable and distributable, even for commercial use. Python has a large community of users and developers, with over two hundred and fifty thousand projects referenced on the Python Package Index (<https://pypi.org>).

Phenotype	PPMI	DIGPD
Depression	GDS	HADS
Anxiety	STAI	HADS
Parkinson's Disease	MDS-UPDRS	MDS-UPDRS
REM Sleep Behavior Disorder	RBDSQ	Binary variable
Cognition	MoCA	MMSE
Activities of Daily Living	Schwab and England ADL	Schwab and England ADL
Autonomic dysfunction	SCOPA-AUT	SCOPA-AUT
Impulsive behaviors	QUIP	Investigator diagnosis
Sleepiness	ESS	ESS
Non-motor symptoms	NMSS	NMSS

Table 1.3: Questionnaires and scales used in PPMI and DIGPD.

ADL: Activities of Daily Living; DIGPD: Drug Interaction With Genes in Parkinson's Disease; ESS: Epworth Sleepiness Scale; GDS: Geriatric Depression Scale; HADS: Hospital Anxiety and Depression Scale; MDS-UPDRS: Movement Disorder Society-sponsored revision of the Unified Parkinson's Disease Rating Scale; MoCA: Montreal Cognitive Assessment; MMSE: Mini-Mental State Examination; NNMS: Non-Motor Symptoms Scale for Parkinson's Disease PPMI: Parkinson's Progressive Markers Initiative; SCOPA-AUT: Scales for Outcomes in Parkinson's Disease – Autonomic Questionnaire; STAI: State-Trait Anxiety Inventory; UPSIT: University of Pennsylvania Smell Identification Test.

A large part of the growing popularity of Python is due to the increasing interest in data science and the availability of maintained, well-documented, high-quality Python packages for science. From data manipulation to machine learning to data visualization, Python is the *de facto* programming language for data science. We will briefly introduce the packages that we used to perform analyses.

NumPy `numpy` (Harris et al., 2020a) is the fundamental package for scientific computing with Python. Fast and versatile, the NumPy vectorization, indexing, and broadcasting concepts are the standards of array computing today. `numpy` offers many numerical computing tools: comprehensive mathematical functions, random number generators, linear algebra routines, Fourier transforms, and more.

SciPy `scipy` (Virtanen et al., 2020) is a package dedicated to scientific computing. It provides many user-friendly and efficient numerical routines, such as routines for numerical integration, interpolation, optimization, linear algebra, and statistics, as well as sparse matrices.

pandas `pandas` (McKinney, 2010) is a fast, powerful, flexible and easy to use open source data analysis and manipulation tool. It provides tools for reading and writing data, fast and efficient data manipulation, and high performance merging and joining of data sets.

Statsmodels `statsmodels` (Seabold and Perktold, 2010) is a Python module that provides classes and functions for the estimation of many different statistical models. It also provides utilities to conduct statistical tests and statistical data exploration.

UMAP `UMAP` (McInnes et al., 2018) is a Python package providing an implementation of the Uniform Manifold Approximation Projection (UMAP) algorithm, which is a popular dimension reduction techniques that can be used for visualization but also for general non-linear dimension reduction (McInnes et al., 2018).

Scikit-learn `scikit-learn` (Pedregosa et al., 2011) is a popular package for machine learning in Python. It is a versatile toolbox for data mining and data analysis, making available numerous machine learning algorithms and utility tools under a unified application programming interface. `scikit-learn` is easily accessible to everybody and usable in various contexts.

XGBoost `xgboost` (Chen and Guestrin, 2016) is an optimized distributed gradient boosting library designed to be highly efficient, flexible and portable. `xgboost` implements machine learning algorithms under the gradient boosting framework that solve many data science problems in a fast and accurate way.

PyTorch `pytorch` (Paszke et al., 2019) is an open source machine learning framework dedicated to deep learning. A rich ecosystem of tools and libraries extends `pytorch` and supports development in computer vision, natural language processing and more.

Matplotlib `matplotlib` (Hunter, 2007) is a comprehensive library for creating static, animated, and interactive visualizations in Python. Several toolkits are available which extend `matplotlib` functionality.

C/C++

C is a general-purpose programming language that is widely used for systems programming in implementing operating systems and embedded system applications. C++, an extension of C, was designed with performance, efficiency, and flexibility as its core.

Thanks to low overhead, C and C++ enable programmers to create efficient implementations of algorithms and data structures, useful for computationally intense programs. Most of the Python packages with intensive computations are partially written in C or one of its variants under the hood. Nonetheless, some scientific libraries are completely written in C/C++, in particular for genetic analyses.

PLINK PLINK (Chang et al., 2015) is a widely used C/C++ tool set for research in population genetics and genome-wide association studies. PLINK provide utilities for data management, basic statistics, linkage disequilibrium calculation, population stratification, association analysis, and tests for epistasis.

GCTA GCTA (Yang et al., 2011) is a C/C++ tool for genome-wide complex trait analysis. GCTA was initially designed to estimate the proportion of phenotypic variance explained by all genome-wide single SNPs for complex traits. It has been subsequently extended for many other analyses to better understand the genetic architecture of complex traits, such as estimation of SNP-based heritability and genomic risk prediction.

R

R is a popular programming language for statistical computing and many packages have been developed to perform statistical analyses. R provides a wide variety of statistical (linear and nonlinear modelling, classical statistical tests, time-series analysis, classification, clustering, etc.) and graphical techniques, and is highly extensible.

meta meta (Balduzzi et al., 2019) is a user-friendly general package providing standard methods for meta-analysis. meta provides fixed effects and random effects models and several plots for meta-analysis among other tools.

Chapter 2

Prediction of impulse control disorders in Parkinson’s disease

This chapter has been submitted to the *Annals of Neurology* journal as:

Johann Faouzi, Samir Bekadar, Baptiste Couvy-Duchesne, Fanny Artaud, Alexis Elbaz, Graziella Mangone, Olivier Colliot and Jean-Christophe Corvol. “Prediction of impulse control disorders in Parkinson’s disease”.

Abstract

Objective To predict the (future) occurrence of impulse control disorders (ICDs) in Parkinson’s disease (PD) using longitudinal data, the first study using cross-validation and replication in an independent cohort.

Methods We included data from two longitudinal PD cohorts (training set: PPMI, Parkinson’s Progression Markers Initiative; replication: DIGPD, Drug Interaction With Genes in Parkinson’s Disease). Patients with at least two visits and with genetic data available were included into the analysis. We trained three logistic regressions and a recurrent neural network to predict ICD at the next visit using clinical risk factors and genetic variants previously associated with ICDs. We quantified performance using the area under the receiver operating characteristic curve (ROC AUC) and average precision.

Results We included 380 PD subjects (2,728 visits) from PPMI and 388 PD subjects (2,101 visits) from DIGPD in our analyses. The number of patients presenting ICDs during follow-up were 143 (38%) in PPMI and 192 (49%) in DIGPD. All the models performed relatively well at predicting the ICDs at the next visit (PPMI: ROC AUC = 0.81 [0.75 - 0.85], DIGPD: ROC AUC = 0.77 [0.67 - 0.80]). Taking previous data from all visits into account improved the predictive performance (PPMI: ROC AUC = 0.83 [0.80 - 0.85], DIGPD: ROC AUC = 0.80 [0.80 - 0.80]), as compared to only using

the baseline visit (PPMI: ROC AUC = 0.75, DIGPD: ROC AUC = 0.67). Recurrent neural networks did not improve the predictive performance.

Interpretation ICDs in PD can be predicted with acceptable accuracy, which may be used to improve the management of PD patients and diminish the potentially devastating impacts of ICDs.

2.1 Introduction

Although Parkinson's disease (PD) is mostly known for its motor symptoms, numerous non-motor symptoms have been reported to occur during the course of the disease (Hiseman and Fackrell, 2017). Impulse control disorders (ICDs), a class of psychiatric disorders characterized by impulsivity, are common in PD, with half of PD cases expected to experience some of them by 5 years after disease onset (Corvol et al., 2018). The four most common ICDs in PD are pathological gambling, compulsive eating, hypersexuality, and compulsive eating disorder. ICDs are associated with reduced quality of life, strained interpersonal relationships, increased caregiver burden, and require prompt addressing (Weintraub and Claassen, 2017). Several case reports suggest that partial and total discontinuations of dopamine agonist (DA) treatment leads to a resolution of ICDs (Mamikonyan et al., 2008; Nirenberg and Waters, 2006).

Many factors have been associated with ICDs in PD, including socio-demographic, clinical and genetic biomarkers (Grall-Bronnec et al., 2018). In particular, men tend to develop more pathological gambling and hypersexuality disorders while women develop more compulsive buying and eating disorders (Weintraub and Claassen, 2017). A younger age has been associated with ICDs in PD in numerous studies (Callesen et al., 2014; Poletti et al., 2013; Pontieri et al., 2015; Weintraub et al., 2010a). Anxiety (Leroi et al., 2012; Pontieri et al., 2015; Voon et al., 2011), depression (Callesen et al., 2014; Voon et al., 2011), and rapid eye movement (REM) sleep behavior disorders (Fantini et al., 2015; Ramírez Gómez et al., 2017) have also been correlated to ICDs. Dopamine replacement therapy, in particular dopamine agonists, has been strongly associated with ICDs. Finally, associations between ICDs and several single-nucleotide polymorphisms (SNPs) in dopamine signaling pathway genes have been suggested (Castro-Martínez et al., 2018; Cormier-Dequaire et al., 2018; Erga et al., 2018; Krishnamoorthy et al., 2016; Lee et al., 2009; Zainal Abidin et al., 2015).

The predictive performance of these factors altogether has been underexplored. Only two studies report predictions at the patient level (Erga et al., 2018; Kraemmer et al., 2016). In both studies, authors trained a logistic regression using clinical and genetic data, and measured its predictive performance using the area under receiver operating characteristic (ROC) curve (ROC AUC). None of these studies had cross-validation or a replication cohort, altering the confidence in the reported performance (Koul et al., 2018).

Our main objective was to predict ICDs from clinical and genetic using machine learning approaches. We utilized two longitudinal cohorts to train and cross-validate the models on one cohort, but also assess the generalization capability of these models on the second cohort. The objective was to predict the risk of ICDs at the next visit, knowing the clinical history of the patient and their genotyping data.

2.2 Materials and methods

2.2.1 Populations

We used data from two research cohorts: the Parkinson’s Progression Markers Initiative (PPMI) database and the Drug Interaction With Genes in Parkinson’s Disease (DIGPD) study.

PPMI (<https://www.ppmi-info.org>) is a multicenter observational clinical study using advanced imaging, biologic sampling, and clinical and behavioral assessments to identify biomarkers of PD progression (Marek et al., 2011). Data was gathered during face-to-face visits every 6-12 months. PD subjects were de novo and drug-naïve at baseline. We downloaded the clinical and genetic data from the PPMI database (<https://www.ppmi-info.org/data>) on the 17th of October, 2019.

DIGPD is a French multicenter longitudinal cohort with annual follow-up of PD patients (Corvol et al., 2018). Eligible criteria consist in recent PD diagnosis (UK Parkinson’s Disease Society Brain Bank criteria) with disease duration less than 5 years at recruitment. Data was gathered during face-to-face visits every 12 months following standard procedures.

Both studies were conducted according to good clinical practice, obtained approval from local ethic committees and regulatory authorities, and all patients provided informed consent prior to inclusion.

2.2.2 Participants and clinical measurements

Inclusion criteria consisted of having: (i) a PD diagnosis, (ii) a baseline visit and at least another visit, (iii) clinical and genetic data available, and (iv) PD medication taken available.

We included socio-demographics and clinical variables that have been associated with ICDs in the literature: age of PD onset, length of follow-up, sex, past ICDs, continuous scales of anxiety, depression and REM sleep, and the motor exam (part III) of the Movement Disorders Society-sponsored revision of the Unified Parkinson’s Disease Rating Scale (MDS-UPDRS). ICDs were assessed at each visit using the Questionnaire for Impulsive-Compulsive Disorders in Parkinson’s Disease - Rating Scale (Weintraub et al., 2012) in PPMI, and through semi-structured interviews by a movement disorder specialist in DIGPD. We standardized each feature since some of them were assessed

with different scales and because it is a common requirement for most machine learning estimators.

We took into account PD medication with three binary variables corresponding to the main classes of treatment (levodopa, dopamine agonists, others) and we derived more specific variables for dopamine agonists: mean daily, maximum daily and total doses (expressed in levodopa equivalent) and cumulative duration.

2.2.3 Genetic variants

In absence of genome-wide association study on ICDs in PD, we considered 50 genetic variants selected as previously described: 20 variants from 16 genes involved in dopamine, serotonin, glutamate, norepinephrine and opioid systems and previously associated with ICD in PD or in the general population (Cormier-Dequaire et al., 2018); 30 additional variants from 10 genes differentially expressed after an acute challenge of levodopa in the striatum in a mouse model of dopamine denervation (Charbonnier-Beaupel et al., 2015).

Genotyping data were collected using NeuroX (Nalls et al., 2015) arrays in PPMI (267,607 variants measured), and Illumina Multi-Ethnic Genotyping Arrays in DIGPD (1,779,819 variants). We excluded variants with missing rates greater than 2% and variants deviating from Hardy-Weinberg equilibrium ($p < 10^{-8}$). We excluded related individuals (third-degree family relationships), individuals with mismatch between reported sex and genetically determined sex, and individuals with outlying heterozygosity (± 3 standard deviation). We imputed missing SNPs using the Michigan Imputation Server (Das et al., 2016) for PPMI and the Sanger Imputation Server (McCarthy et al., 2016) for DIGPD, using the reference panel of the Haplotype Reference Consortium (release 1.1) (McCarthy et al., 2016). We filtered variants based on their imputation quality ($R^2 > 0.6$ for PPMI, INFO score > 0.9 for DIGPD).

2.2.4 Data processing

Processing genetic data and extracting variants of interest matching inclusion criteria was performed using the PLINK (Chang et al., 2015) software. Processing of the different text-like files was performed using the `pandas` (McKinney, 2010) and `NumPy` (Harris et al., 2020a) Python packages. Missing values were imputed in a forward-fill fashion: for a given subject and a given feature, missing values were imputed using the most recent non-missing value for this subject and this feature. Baseline missing values were imputed using the mean baseline values on the training set.

2.2.5 Machine learning algorithms

We investigated five standard machine learning algorithms implemented in the `scikit-learn` (Pedregosa et al., 2011) and `XGBoost` (Chen and Guestrin, 2016) Python packages:

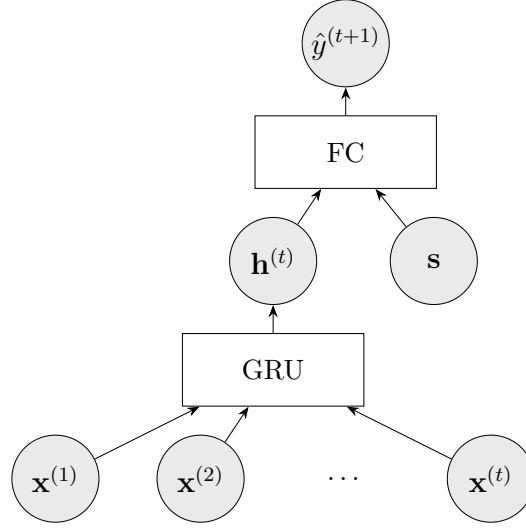


Figure 2.1: Architecture of the recurrent neural network. The clinical features assessed at several visits are used as input of the Gated Recurrent Unit (GRU). The GRU extracts information from these clinical features into a vector $\mathbf{h}^{(t)}$. This vector and the time-independent variables, namely the socio-demographic and genetic data denoted as \mathbf{s} , are used as input of a Fully Connected (FC) function followed by a sigmoid activation, returning the probability of having an impulse control disorder at the next visit.

logistic regression, support vector machines with a linear kernel and a RBF kernel (Boser et al., 1992; Cortes and Vapnik, 1995), random forest (Breiman, 2001) and gradient tree boosting (Friedman, 2001; Mason et al., 2000). These algorithms expect a fixed number of features as input. In order to deal with varying numbers of visits, we reduced all the previous visits into one “summary” visit using a convex combination. A convex combination is a linear combination such that the weights are all non-negative and sum to one. The weights indicate how much each visit contributes to this “summary” visit. A weight of 1 for the first visit means that the “summary” visit is simply the baseline visit, while a weight of 1 for the latest visit means that the summary visit is simply the most recent visit. One can also give uniform weights, so that each visit contributes equally to this summary visit, or higher weights to most recent visits if they are assumed to be more important than older visits.

As the prediction task is longitudinal, we also investigated the use of recurrent neural networks. Recurrent neural networks are a class of artificial neural networks dedicated to sequential data. We employed a simple architecture (Figure 2.1) with a Gated Recurrent Unit (Cho et al., 2014) to extract information from the clinical measurements, followed by a concatenation of this vector with the socio-demographic and genetic data, followed by a Fully Connected function with a sigmoid activation. We use the PyTorch (Paszke et al., 2019) Python package to build and train the recurrent neural network.

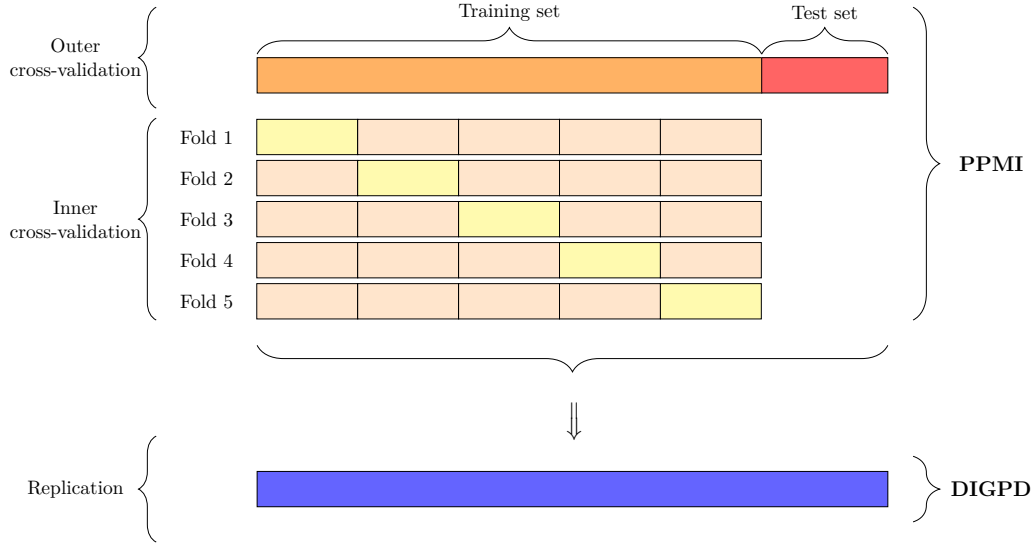


Figure 2.2: Cross-validation procedure. We employed a nested cross-validation procedure on the PPMI cohort. In the outer loop, we split the PPMI subjects into training and test sets, while the inner loop was a 5-fold subject-level cross-validation to optimize the hyper-parameters of the model. The model with the optimal values for the hyper-parameters was evaluated on the test set of PPMI and on the whole DIGPD cohort.

2.2.6 Cross-validation

We used PPMI as the training (discovery) cohort, and DIGPD as the testing (replication) cohort. To unbiasedly estimate the predictive performance of the models, we employed a nested cross-validation procedure that is illustrated in [Figure 2.2](#). In the outer loop, we randomly split 80% of the PPMI subjects into the training set and the remaining 20% into the test set. In the inner loop, we performed a 5-fold subject-level cross-validation procedure to optimize the hyper-parameters of the models on the training set. These hyper-parameters control how the algorithms fit the training data. For instance, these hyper-parameters included the type (l1 or l2 penalty) and amount (lambda parameter) of regularization for the linear models. In particular, logistic regression models were regularized. After finding the optimal values for the hyper-parameters, each model was evaluated on the test set. Finally, we evaluated the performance of each model on the whole DIGPD cohort.

2.2.7 Statistical analysis

Baseline characteristics in both cohorts were compared with chi-squared tests for categorical variables and t-tests for continuous variables using the `SciPy` ([Virtanen et al., 2020](#)) Python package. Predictive performance was mainly evaluated using the area under the receiver operating characteristic curve and average precision (AP). AP summarizes a precision-recall curve as the weighted mean of precisions achieved at each

threshold, with the increase in recall from the previous threshold used as the weight. The precision-recall curve is similar to the ROC curve, but plots the precision (positive predictive value) against the recall (sensitivity). The precision-recall curve does not take into account the true negatives, and is particularly useful when the positives are more important than the negatives (false negatives are more serious than false positives). Other metrics included accuracy, balanced accuracy, sensitivity and specificity. ROC and precision-recall curves were plotted using the `matplotlib` (Hunter, 2007) Python package, and all the metrics were computed using the `scikit-learn` (Pedregosa et al., 2011) package. Comparison between ROC AUC was measured using DeLong test (DeLong et al., 1988). P-values were adjusted for multiple comparisons using Bonferroni correction. We did not statistically compare AP scores as we were not aware of a relevant statistical test to do so, and reported 95% confidence intervals using bootstrapping.

2.3 Results

2.3.1 Population characteristics

Out of the 423 PD subjects in PPMI, we excluded 1 subject for not having a baseline visit, 2 for not having medication records and 40 for not having genetic data. Out of the 415 PD subjects in DIGPD, we excluded 27 for having only a baseline visit. No subjects were excluded based on their genetic data. Thus, we included 380 PD subjects from PPMI and 388 PD subjects from DIGPD in our analyses. The 380 PPMI subjects had a total of 2,728 visits, while the 388 DIGPD subjects had a total of 2,101 visits. Since our objective was to predict the occurrence of ICDs at the next visit, the number of observations for a given subject is equal to their number of visits minus 1. Thus, the total number of observations was equal to 2,348 in PPMI and 1,713 in DIGPD.

Clinical characteristics are presented in Table 2.1. Age and sex in both cohorts were not significantly different. PPMI subjects had significantly more visits and smaller intervals between back-to-back visits, as well as longer follow-ups. DIGPD subjects had significantly lower scores in the motor exam of the MDS-UPDRS. The prevalence of ICDs at baseline was significantly higher in DIGPD than in PPMI, as well as their lifetime prevalence. Both differences might be explained by the fact that PD subjects are de novo and drug-naïve at baseline in PPMI whereas they are not in DIGPD. Other phenotypes (anxiety, depression, and REM sleep disorders) were not statistically compared due to the different scales used.

Concerning genetic data, we excluded 1 genetic variant for being a variable number of tandem repeat polymorphism. Furthermore, we excluded 18 SNPs for having too low imputation quality scores. Finally, 31 SNPs were included in our analyses (Table A.1).

Characteristic	PPMI	DIGPD	<i>p</i> -value
Age (in years)	60.67 \pm 9.71	58.99 \pm 9.75	1.71×10^{-2}
Sex (F/M)	127/253 (33%)	155/233 (40%)	7.16×10^{-2}
Length of follow-up (in years)	5.86 \pm 1.95	4.82 \pm 1.83	7.07×10^{-14}
Number of visits per subject	7.18 \pm 2.96	5.41 \pm 1.66	4.21×10^{-35}
Interval between visits (in years)	0.95 \pm 0.35	1.090.33	3.24×10^{-39}
Anxiety	STAI: 93.55 \pm 7.96	HAD: 6.82 \pm 3.77	
Depression	GDS: 5.25 \pm 1.47	HAD: 4.59 \pm 3.16	
REM sleep	RBDSQ: 4.17 \pm 2.71	1/0: 86/302(22%)	
MDS-UPDRS III	20.87 \pm 8.86	9.91 \pm 5.33	3.77×10^{-72}
Baseline ICD (1/0)	42/338 (13%)	76/312 (20%)	5.34×10^{-3}
Lifetime ICD (1/0)	143/237 (38%)	192/196 (49%)	1.12×10^{-3}

Table 2.1: Baseline characteristics. For continuous variables, mean \pm standard deviation is reported. For binary variables, the count for both categories is reported as well as the proportion of the first category. Statistical differences were assessed using independent t tests for continuous variables and chi-squared tests for categorical variables.

HAD: Hospital Anxiety and Depression Scale; ICD: Impulse control disorders; MDS-UPDRS: Movement Disorders Society-sponsored revision of the Unified Parkinson's Disease Rating Scale; QUIP: Questionnaire for Impulse-Compulsive Disorders in Parkinson's Disease; REM: Rapid eye movement; RBDSQ: Rapid eye movement Sleep Behavior Disorder Screening Questionnaire; STAI: State-Trait Anxiety Inventory.

2.3.2 Predictive performance

Table 2.2 presents the predictive performance for the four main models: logistic regression using the baseline visit, the most recent visit, and the mean over all the past visits, and the recurrent neural network. The logistic regression using the baseline visit had the lowest scores on both cohorts (ROC AUC = 0.75 and AP = 0.44 in PPMI, ROC AUC = 0.67 and AP = 0.43 in DIGPD). By contrast, the recurrent neural network yielded the highest scores in PPMI (ROC AUC = 0.85, AP = 0.61), while the logistic regression using the most recent visit yielded the highest scores in DIGPD (ROC AUC = 0.802, AP = 0.64). Figure 2.3 and Figure 2.4 show the ROC and precision-recall curves for the four main models in PPMI and in DIGPD respectively. The recurrent neural network models had sensitivities of 61% and 70% and specificities of 90% and 82% in PPMI and DIGPD respectively, at the default threshold (probability > 0.5).

		Logistic regression using only the baseline visit	Logistic regression using only the previous visit	Logistic regression using the mean over past visits	Recurrent neural network
ROC AUC	PPMI	0.75 ([0.69, 0.81])	0.80 ([0.73, 0.86])	0.84 ([0.77, 0.89])	0.85 ([0.79, 0.90])
	DIGPD	0.67 ([0.64, 0.70])	0.80 ([0.78, 0.83])	0.80 ([0.77, 0.82])	0.80 ([0.78, 0.83])
Average precision	PPMI	0.44 ([0.33, 0.56])	0.45 ([0.36, 0.58])	0.60 ([0.49, 0.72])	0.61 ([0.49, 0.73])
	DIGPD	0.43 ([0.39, 0.48])	0.64 ([0.60, 0.69])	0.62 ([0.57, 0.67])	0.62 ([0.58, 0.68])
Accuracy	PPMI	0.77 ([0.74, 0.81])	0.82 ([0.78, 0.85])	0.84 ([0.80, 0.87])	0.86 ([0.83, 0.89])
	DIGPD	0.57 ([0.54, 0.59])	0.59 ([0.57, 0.61])	0.64 ([0.62, 0.67])	0.78 ([0.76, 0.80])
Balanced accuracy	PPMI	0.69 ([0.63, 0.75])	0.76 ([0.71, 0.82])	0.77 ([0.70, 0.82])	0.76 ([0.70, 0.82])
	DIGPD	0.60 ([0.58, 0.63])	0.67 ([0.65, 0.69])	0.69 ([0.67, 0.72])	0.76 ([0.73, 0.78])
Sensitivity	PPMI	0.57 ([0.46, 0.68])	0.69 ([0.58, 0.79])	0.66 ([0.54, 0.77])	0.61 ([0.50, 0.73])
	DIGPD	0.68 ([0.64, 0.73])	0.84 ([0.80, 0.87])	0.81 ([0.77, 0.84])	0.70 ([0.66, 0.74])
Specificity	PPMI	0.81 ([0.77, 0.85])	0.84 ([0.81, 0.88])	0.87 ([0.84, 0.91])	0.90 ([0.87, 0.93])
	DIGPD	0.52 ([0.49, 0.55])	0.50 ([0.47, 0.53])	0.58 ([0.56, 0.61])	0.82 ([0.80, 0.84])

Table 2.2: Results of the four main models. Predictive performance for the four main models on both cohorts are reported. 95% confidence intervals were estimated using 2000 bootstrap samples.

DIGPD: Drug Interaction With Genes in Parkinson’s Disease; PPMI: Parkinson’s Progression Markers Initiative; ROC AUC: area under the receiver operating characteristic curve.

All four models were statistically better than random guess ($p < 0.001$). Logistic regression using only the baseline visit was statistically worse than at least two other models in both cohorts (Figure 2.5). The three other models were not statistically different from each other in both cohorts.

Although AP scores for the three best models were higher in DIGPD than in PPMI, the prevalence of ICDs, computed over all the (patient, visit) pairs, was twice higher in DIGPD than in PPMI (27% in DIGPD, 14% in PPMI). As AP scores of random guess are equal to the prevalence of the positive class, the differences between AP scores in both cohorts should be interpreted with much caution.

The other machine learning algorithms (support vector machines with linear and RBF kernels, random forest, and gradient tree boosting) and other reduction approaches (giving positive weights to all the past visits, but higher weights to more recent visits) yielded comparable results (Table A.3 and Table A.4).

To evaluate the impact of the splitting of PPMI into training and test sets on the predictive performance, we repeated the cross-validation procedure 10 times and also evaluated the 10 models on DIGPD. All iterations yielded comparable results (Table A.5 and Table A.6).

2.3.3 Contribution of the different features

Since the genetic factors of ICDs in PD are mostly unknown and genotyping data is not usually collected in clinical routine, we investigated the predictive performance of the same algorithms without the genetic variants as input, in order to assess their added value in the models. Table 2.3 presents the ROC AUC of the models with and without genetic variants and their statistical comparison. Only one comparison was statistically different: the logistic regression model using the most recent visit had a higher ROC AUC with genetic variants than without genetic variants (ROC AUC = 0.80 with genetic variants, ROC AUC = 0.79 without genetic variants, $p < 0.001$). The genetic variants did not seem to be major contributors to the decision function of the logistic regression models.

Table 2.4 presents the coefficients of the three logistic regression models without genetic variants as input (see Table A.2 for the coefficients of the three logistic regression models with genetic variants as input). As the logistic regression model using the baseline visit performed significantly worse, and the variables for PD medication were all null (PD patients in PPMI are de novo drug-naïve at baseline, and the medical history of PD patients in DIGPD was not available before their baseline visit), we only interpreted the other two models. The following features had positive coefficients: sex, past ICDs, depression, REM sleep, motor exam, being on other PD medication than levodopa and dopamine agonists, and maximum dose and cumulative duration of dopamine agonists. On the other hand, the following features had negative coefficients:

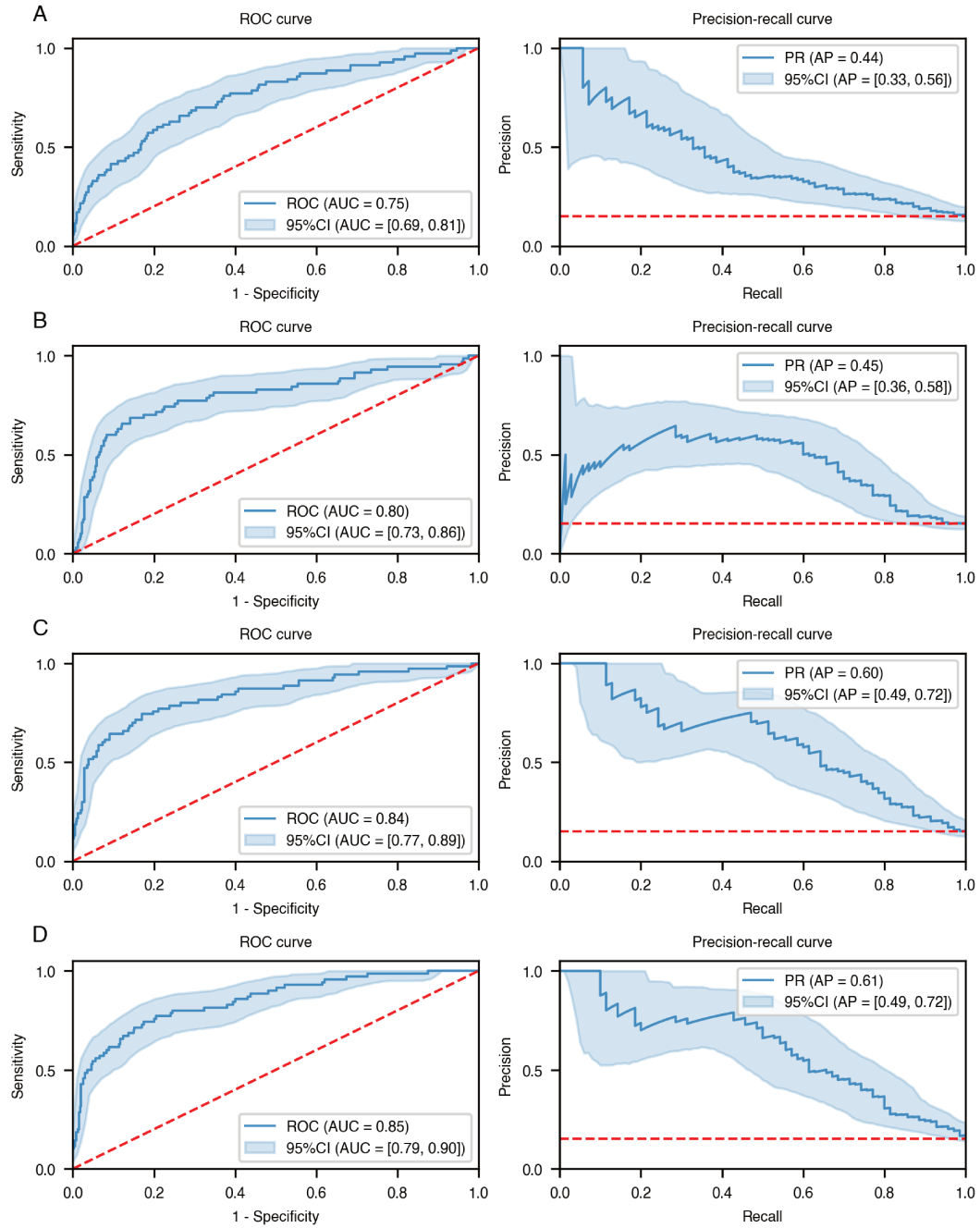


Figure 2.3: ROC and Precision-recall curves on PPMI. The ROC and precision-recall curves were computed for the four models, with 95% confidence intervals behind estimated using 2000 bootstrapping samples: (A) logistic regression using the baseline visit; (B) logistic regression using the most recent visit; (C) logistic regression using the mean over the past visits; (D) recurrent neural networks using all the past visits.

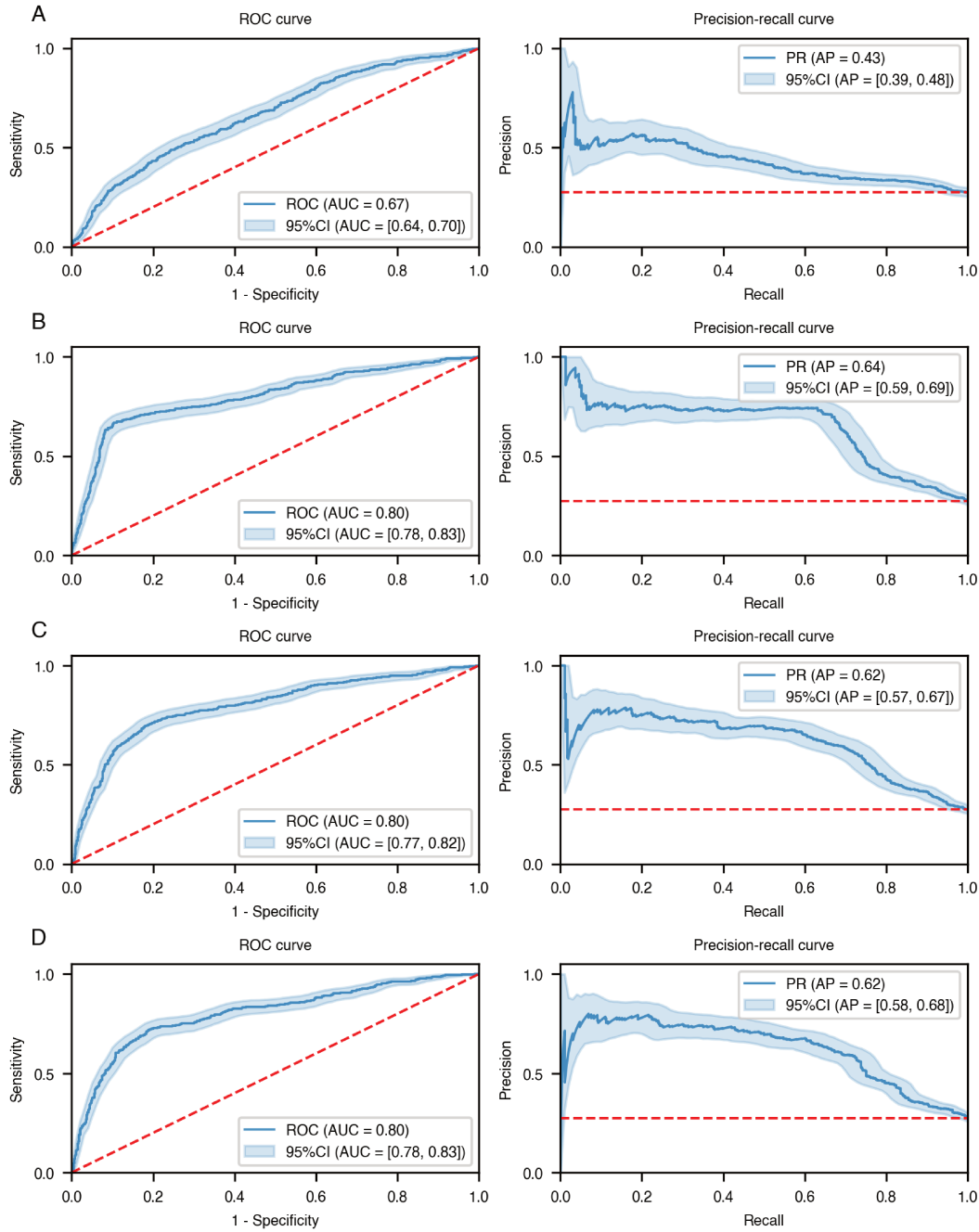


Figure 2.4: ROC and Precision-recall curves on DIGPD. The ROC and precision-recall curves were computed for the four models, with 95% confidence intervals behind estimated using 2000 bootstrapping samples: (A) logistic regression using the baseline visit; (B) logistic regression using the most recent visit; (C) logistic regression using the mean over the past visits; (D) recurrent neural networks using all the past visits.

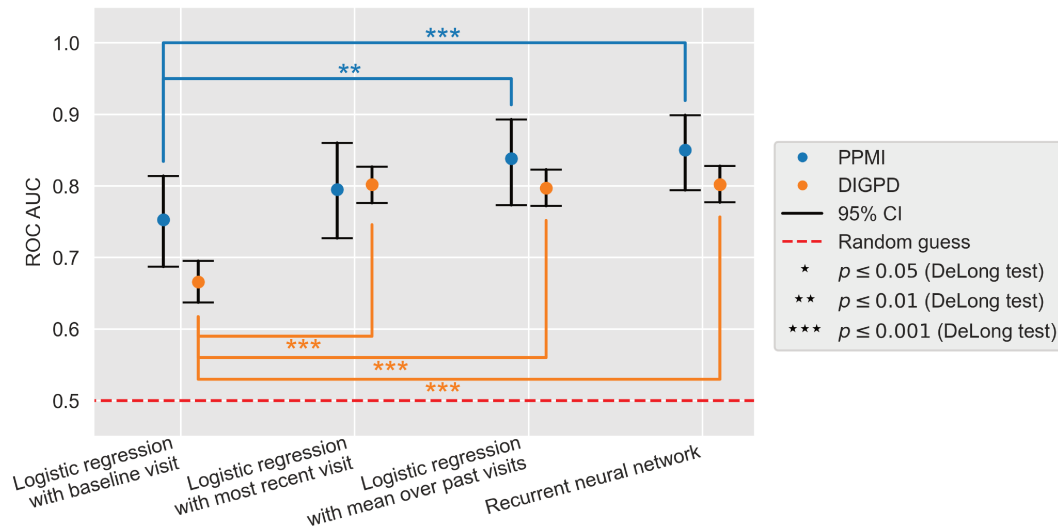


Figure 2.5: Statistical comparison of ROC AUC for the four main models. 95% confidence intervals were computed using 2000 bootstrap samples. P-values were computed using the DeLong test. P-values below the 0.05 threshold after adjustment for multiple comparison using Bonferroni correction are highlighted with at least one asterisk: $p \leq 0.05$ (*), $p \leq 0.01$ (**), $p \leq 0.001$ (***).

CI: confidence interval; DIGPD: Drug Interaction With Genes in Parkinson's Disease; PPMI: Parkinson's Progression Markers Initiative; ROC AUC: area under the receiver operating characteristic curve.

		Logistic regression using only the baseline visit	Logistic regression using only the previous visit	Logistic regression using the mean over past visits	Recurrent neural network
PPMI	With genetic variants	0.753 ([0.687, 0.814])	0.795 ([0.727, 0.860])	0.838 ([0.773, 0.893])	0.850 ([0.794, 0.899])
	Without genetic variants	0.751 ([0.683, 0.816])	0.807 ([0.745, 0.865])	0.843 ([0.782, 0.894])	0.845 ([0.789, 0.895])
	<i>p</i> -value	0.885	0.0971	0.637	0.451
DIGPD	With genetic variants	0.666 ([0.637, 0.695])	0.802 ([0.776, 0.827])	0.797 ([0.772, 0.823])	0.802 ([0.777, 0.828])
	Without genetic variants	0.682 ([0.653, 0.710])	0.786 ([0.758, 0.812])	0.788 ([0.763, 0.813])	0.803 ([0.778, 0.828])
	<i>p</i> -value	0.0143	0.000252	0.109	0.757

Table 2.3: Statistical comparison of ROC AUC for the four main models with and without genetic variants. Differences in ROC AUC between the models with and without genetic variants were assessed with the DeLong test. Significant differences after Bonferroni correction are highlighted in bold font.

DIGPD: Drug Interaction With Genes in Parkinson’s Disease; PPMI: Parkinson’s Progression Markers Initiative; ROC AUC: area under the receiver operating characteristic curve.

		Baseline visit	Most recent visit	Mean over past visits
Socio-demographic	Sex	0.174	0.182	0.344
	Age	-0.318	-0.312	-0.184
Clinical	Past ICDs	1.582	1.921	3.958
	Depression	0.000	0.131	0.887
	State anxiety	0.000	-0.358	0.000
	Trait anxiety	0.000	-0.332	-0.562
	REM sleep	0.723	0.533	0.639
	Motor exam	0.000	0.128	0.343
PD medication	On levodopa	0.000	-0.029	-0.300
	On dopamine agonists	0.000	0.063	0.000
	On other PD medication	0.000	0.171	0.063
	Mean daily dose of dopamine agonists	0.000	-0.130	-0.049
	Maximum daily dose of dopamine agonists	0.000	0.253	0.563
	Total dose of dopamine agonists	0.000	-0.393	-3.392
	Cumulative duration on dopamine agonists	0.000	0.355	2.109
Time to prediction	Time to prediction	0.106	0.032	0.039

Table 2.4: Coefficients of the three logistic regression models without genetic variants as input.

age, anxiety, being on levodopa, and mean daily and total dose of dopamine agonists. The features corresponding to being on dopamine agonists and time to prediction had coefficients close to zero. The variables with the largest absolute values were past ICDs, and total dose and cumulative duration of dopamine agonists.

2.4 Discussion

To the best of our knowledge, this study is the first one evaluating the predictability of ICDs in PD in an unbiased manner using two longitudinal cohorts, including one independent replication cohort.

Two previous studies reported ROC AUC for a prediction task of ICDs in PD (Erga et al., 2018; Kraemmer et al., 2016). Kraemmer and colleagues reported ROC AUC of 0.65 (95% CI 0.58-0.73) with clinical variables only and of 0.76 (95% CI 0.70-0.83) with clinical and genetic variables, while Erga and colleagues reported ROC AUC of 0.68 (95% CI 0.59-0.78) with clinical features only and of 0.70 (95% CI 0.61-0.79) with clinical and genetic features. However, the methods and prediction tasks were different. Erga and colleagues performed a cross-sectional analysis of 119 PD patients from the Norwegian ParkWest study, while Kraemmer and colleagues performed a longitudinal analysis of 276 PD patients from PPMI. In their studies, each patient corresponds to a unique observation, leading to much lower sample sizes. Moreover, both studies did not use cross-validation and did not have a replication cohort, which might lead to overly optimistic reported results (Koul et al., 2018). By contrast, our study uses cross-validation on the training cohort (PPMI) and has a replication cohort (DIGPD), with comparable results in both cohorts. Moreover, both cohorts have different characteristics (de novo drug-naïve patients in PPMI, already-treated patients in DIGPD) and some variables (anxiety, depression, REM sleep) were not measured with the same instruments, suggesting good generalizability of the models in different settings. We also detail our methodology and provide the coefficients of the logistic regression models, so that others can investigate the predictive performance of our models in their cohorts.

The logistic regression coefficients were overall consistent with the literature. For the socio-demographic variables, sex and age have respectively a positive and negative coefficients, in accordance with a younger age and a male sex previously associated with ICDs in PD (Weintraub and Claassen, 2017). Depression, REM sleep and motor exam scores also had positive coefficients, consistent with their positive association (Grall-Bronnec et al., 2018). Anxiety scores had negative coefficients although previously reported to be positively associated with ICD. The maximum dose and cumulative duration of dopamine agonists had positive coefficients, confirming the important role of the dose and the duration of dopamine agonist therapy in the risk to develop ICDs in PD (Corvol et al., 2018). Interestingly, the types of PD medication sparsely contributed to the decision function of the models, with very small coefficients. The mean daily and

total doses of dopamine agonists had negative coefficients, although these coefficients were almost null for the mean daily dose. These derived features have rarely been investigated altogether, making the comparison with the literature difficult. It should be noted that the coefficients are estimated altogether and that the logistic regression models were regularized, so interpretation should be performed with caution.

Although the predictive performance of the models can be considered acceptable, the use of such models in clinical practice would deserve improving their accuracy. We used features that have been associated with ICDs in PD as input of our models, but there are probably more unknown risk factors to be discovered. In addition, as a class of psychiatric disorders, ICDs are particularly complex, with qualitative environmental factors that might play important roles, are difficult to measure, and are not captured by clinical scales used in PD. Assessment of ICDs may also be noisy (e.g. patients hiding or not aware of their behavior), and thus ICDs are probably less predictable in practice than other comorbidities in PD, such as dementia (Liu et al., 2017). Finally, little is known about the genetic factors of ICDs in PD. In absence of genome-wide association study and genetic risk scores for ICDs in PD, we used associated genetic variants from candidate gene analyses (Cormier-Dequaire et al., 2018; Erga et al., 2018; Kraemmer et al., 2016). As variation in complex traits is caused by numerous genetic variants, such analyses have important limitations and many association studies could not be replicated, particularly in psychiatric conditions like schizophrenia (Johnson et al., 2017). More studies, in particular genome-wide association studies, are needed to better understand the genetic landscape of ICDs in PD.

We used ROC AUC as the main metric to evaluate the models, and recurrent neural network models did not have much added value over logistic regression models for this metric. ROC AUC is the area under the ROC curve, plotting the sensitivity against the specificity, and summarizes how much sensitivity and specificity change for different thresholds. However, in practice, a single threshold is generally used. Using the default threshold (probability > 0.5) yielded higher balanced accuracy (mean of sensitivity and specificity) scores in the replication cohort for the recurrent neural network model than the logistic regression models, whereas this was not observed in the training cohort. This might suggest better generalizability of the recurrent neural network model than the logistic regression models when using a single threshold.

Being able to predict ICDs is of critical importance due to their potential medical, financial, and/or legal medical complications. Identifying patients at high risk to develop ICDs at the next visit may lead to changes in the dopaminergic treatment strategy (e.g. decrease the dose of dopamine agonists and increase levodopa) and/or recommend a closer monitoring of behavioral changes by the caregiver. The efficacy of such preventive strategies based on a predictive model remains however to be evaluated. In this perspective, the model may be adapted depending on the relative importance for identifying positives (patients who will develop ICDs) or negatives (patients who will not

develop ICDs). The balanced accuracy scores were equal to 76% for the recurrent neural network models in both cohorts, but the sensitivities (61% vs 70%) and specificities (90% vs 82%) differed, which might be explained by the different prevalences in both cohorts. Using the default threshold (probability > 0.5) made the models more specific than sensitive, which might be a limitation if finding the positives is more important than the negatives. On the other hand, models being more specific than sensitive might be more relevant if the main objective is to propose treatment changes only to patients who are really at risk, and avoid unnecessary modifications in more patients. The threshold can still be adjusted depending on the main objective. Prospective studies are required to validate the models and allow their relevance in clinical routine.

Our study has several limitations. First, the sample sizes are relatively small, in particular on the test set of PPMI due to the use of cross-validation, leading to large confidence intervals. Second, each observation is a (subject, visit) pair and thus the observations are not independent (the intra-subject observations are not independent, but the inter-subject observations are independent), which could lead to underestimating p-values when assessing the statistical difference between ROC AUC. Third, in absence of genome-wide association study and genetic risk scores for ICDs in PD, we used associated genetic variants from candidate gene analyses. Genetic risk scores are more robust estimators of the genetic liability of a phenotype and should be preferred when available ([Wray et al., 2007](#)).

In conclusion, our study shows that ICDs in PD can be predicted with a relatively good accuracy. The developed models were unbiasedly evaluated in two research cohorts, with comparable results. Our study highlights the utility of machine learning to automatically extract information from data and its potential to improve patient care.

Chapter 3

Exploratory analysis of the genetics of impulse control disorders in Parkinson’s disease using genetic risk scores

This chapter has been submitted to the *Parkinsonism and Related Disorders* journal as:

Johann Faouzi, Baptiste Couvy-Duchesne, Samir Bekadar, Olivier Colliot and Jean-Christophe Corvol. “Exploratory analysis of the genetics of impulse control disorders in Parkinson’s disease using genetic risk scores”.

Abstract

Objective To study the association between impulse control disorders (ICDs) in Parkinson’s disease (PD) and genetic risk scores (GRS) for 40 known or putative risk factors (e.g. depression, personality traits).

Background In absence of published genome-wide association studies (GWAS), little is known about the genetics of ICDs in PD. GRS of related phenotypes, for which large GWAS are available, may help shed light on the genetic contributors of ICDs in PD.

Methods We searched for GWAS on European ancestry populations with summary statistics publicly available for a broad range of phenotypes, including other psychiatric disorders, personality traits, and simple phenotypes. We separately tested their predictive ability in two of the largest PD cohorts with clinical and genetic available: the Parkinson’s Progression Markers Initiative database (N = 368, 33% female, age range = [33 - 84]) and the Drug Interaction With Genes in Parkinson’s Disease study (N=373, 40% female, age range = [29 - 85]).

Results We considered 40 known or putative risk factors for ICDs in PD for which large GWAS had been published. After Bonferroni correction for multiple comparisons, no GRS or the combination of the 40 GRS were significantly associated with ICDs from the analyses in each cohort separately and from the meta-analysis.

Conclusion Albeit unsuccessful, our approach will gain power in the coming years with increasing availability of genotypes in clinical cohorts of PD, but also from future increase in GWAS sample sizes of the phenotypes we considered. Our approach may be applied to other complex disorders, for which GWAS are not available or limited.

3.1 Introduction

Although the cardinal symptoms of Parkinson's disease (PD) are motor, many non-motor symptoms frequently occur during the course of the disease, including psychiatric comorbidities. Impulse control disorders (ICDs), a class of psychiatric disorders characterized by impulsivity, are common in PD, with half of PD cases expected to experience them within 5 years of the disease onset (Corvol et al., 2018). The four most common ICDs in PD are pathological gambling, compulsive eating, hypersexuality, and compulsive eating disorder. ICDs are associated with reduced quality of life, strained interpersonal relationships, increased caregiver burden, and require prompt management (Weintraub and Claassen, 2017).

Numerous factors have been associated with ICDs in PD, including socio-demographic, clinical and genetic variables (Grall-Bronnec et al., 2018). Associations from candidate gene analyses between ICDs and several genetic variants have been reported in the following genes: *ANKK1* (Hoenicka et al., 2015), *DAT1* (Cormier-Dequaire et al., 2018), *DRD1* (Erga et al., 2018; Zainal Abidin et al., 2015), *DRD2* (Kraemmer et al., 2016; Zainal Abidin et al., 2015), *DRD3* (Castro-Martínez et al., 2018; Krishnamoorthy et al., 2016; Lee et al., 2009), *GRIN2B* (Lee et al., 2009; Zainal Abidin et al., 2015), *HTR2A* (Kraemmer et al., 2016; Lee et al., 2012), *OPRK1* (Cormier-Dequaire et al., 2018; Kraemmer et al., 2016), *OPRM1* (Cormier-Dequaire et al., 2018), and *SLC22A1* (Redenek et al., 2019). Several studies also reported no consistent associations with variants from some of the same genes (Cormier-Dequaire et al., 2018; Vallelunga et al., 2012), highlighting the variability and the lack of replication of the reported associations.

Variation in complex traits is caused by numerous genetic variants. Each genetic variant usually provides limited information because the relative causal risk of each variant is small (Wray et al., 2007). On the other hand, the combined risk of numerous low-risk variants can explain a significant proportion of the genetic variance. Genetic risk scores (GRS), obtained from genome-wide association studies (GWAS), linearly summarize the contribution of these numerous variants into a single score. Using such GRS allows for studying traits not collected in the PD cohorts or diseases that would be too rare to allow direct evaluation of the comorbidities.

ICDs in PD have been rarely studied using GRS. Only one study looked for associations between GRS and ICDs in PD (Ihle et al., 2020). In this study, authors computed a GRS of PD using 90 SNPs reaching genome-wide significance in a meta-analysis of 17 GWAS (Nalls et al., 2019) and did not find an association between this GRS and ICDs in PD. Their power was limited due to the small sample size. Furthermore, the GRS may only capture part of the genetic risk factors and would benefit from larger GWAS (Dudbridge, 2013). ICDs are not associated with PD itself (de la Riva et al., 2014), and may be associated with personality traits or psychiatric endophenotypes, which has been little studied.

Our main objective was to evaluate the predictive accuracy of a broad range of GRS in order to shed light on the genetic determinants of ICDs in PD. We were particularly interested in GRS for other psychiatric disorders, but also personality traits, including impulsivity, some of which have been associated with ICDs in PD (Callesen et al., 2014; Sáez-Francàs et al., 2016; Voon et al., 2011).

3.2 Materials and methods

3.2.1 Populations

We used data from two research cohorts: the Parkinson’s Progression Markers Initiative (PPMI) database and the Drug Interaction With Genes in Parkinson’s Disease (DIGPD) study.

PPMI (<https://www.ppmi-info.org>) is a multicenter observational clinical study using advanced imaging, biologic sampling and clinical and behavioral assessments to identify biomarkers of PD progression. Data was gathered during face-to-face visits every 6–12 months. PD subjects were de-novo and drug-naïve at baseline. We downloaded the clinical and genetic data from the PPMI database (<https://www.ppmi-info.org/data>) on the 17th of October, 2019.

DIGPD is a French multicenter longitudinal cohort with annual follow-up of PD patients (Corvol et al., 2018). Eligible criteria consist in recent PD diagnosis (UK Parkinson’s Disease Society Brain Bank criteria) with disease duration less than 5 years at recruitment. Data was gathered during face-to-face visits every 12 months following standard procedures.

Both studies were conducted according to good clinical practice, obtained approval from local ethic committees and regulatory authorities, and all patients provided informed consent prior to inclusion.

3.2.2 Participants

Inclusion criteria in our analyses included having: (i) a PD diagnosis, (ii) at least two visits measuring ICDs, (iii) clinical and genetic data available, and (iv) a European

genetic ancestry. We identified 378 subjects in PPMI and 382 subjects in DIGPD matching the first three criteria.

ICDs were assessed at each visit using the Questionnaire for Impulsive-Compulsive Disorders in Parkinson's Disease - Rating Scale (Weintraub et al., 2012) in PPMI, and through semi-structured interviews by a movement disorders specialist in DIGPD. The ICD phenotype was defined as the lifetime presence of ICDs.

3.2.3 Genetic ancestry

To date, most GWAS have been conducted in populations of European ancestry, which limits the use of GWAS-derived GRS in non-European ancestry populations (Wang et al., 2020), and their transferability to other populations depends on many factors such as linkage disequilibrium, allele frequencies, and genetic architecture. Directly computing GRS in another ancestry group than the one from the corresponding GWAS can lead to biased GRS (Martin et al., 2017).

To estimate the genetic ancestry of the PD subjects in PPMI and DIGPD, we used data from the 1000 genomes (1000G) project to learn a low-dimensional representation of the genetic data, which captures the main dimension of ancestry. Using the 50,842 common raw SNPs between 1000G, PPMI and DIGPD, we applied the Uniform Manifold Approximation Projection (McInnes et al., 2018) (UMAP) algorithm on the 1000G data to learn a low-dimensional space of the raw SNPs. Finally, we projected the PPMI and DIGPD subjects onto the main principal components to identify in which clusters they were the closest to. Subjects projected on another cluster than the European cluster were excluded.

3.2.4 Genotyping and quality control

Genotype data was acquired using NeuroX (Nalls et al., 2015) arrays in PPMI (267,607 variants measured), and Illumina Multi-Ethnic Genotyping Arrays in DIGPD (1,779,819 variants). We excluded variants with missing rates greater than 2% and variants deviating from Hardy-Weinberg equilibrium ($p < 10^{-8}$). We excluded related individuals (third-degree family relationships), individuals with mismatching between reported sex and genetically determined sex, and individuals with outlying heterozygosity (± 3 standard deviation). We then imputed missing SNPs using the Michigan Imputation Server (Das et al., 2016) for PPMI and the Sanger Imputation Server (McCarthy et al., 2016) for DIGPD, using the reference panel of the Haplotype Reference Consortium (release 1.1) (McCarthy et al., 2016).

For GRS calculation, we selected SNPs that were (i) biallelic, (ii) frequent enough (minor allele frequency $> 1\%$), and (iii) imputed with sufficient accuracy ($R^2 > 0.8$ for PPMI, INFO Score > 0.9 for DIGPD).

3.2.5 Phenotypes and genome-wide association studies

Phenotypes of interest included known or putative factors associated with ICDs in PD, such as anxiety, depression, personality traits including impulsivity, eating and sleep disorders. We were also interested in more general phenotypes such as body height, body mass index (BMI), intelligence, and number of years of education, more because of the sample size of the corresponding GWAS rather than their prior association with ICDs in PD. In particular, body height and body mass index are phenotypes that are easy to collect with precision, and for which very large GWAS are available and the corresponding GRS explain a large part of the variance. These phenotypes are also usually collected in research cohorts, allowing for comparing the GRS with the true phenotypes, and thus validating our computation of the GRS.

We used the NHGRI-EBI GWAS Catalog ([Buniello et al., 2019](#)) to select the largest GWAS to date on samples of European ancestry. When summary statistics from several GWAS were available for a given phenotype, we only included the largest study.

3.2.6 Computation of genetic risk scores

When summary statistics were fully available, we estimated the coefficients of the GRS using the SBLUP ([Robinson et al., 2017](#)) algorithm implemented in the GCTA ([Yang et al., 2011](#)) software. SBLUP directly estimates GRS coefficients from summary statistics, using a reference sample to estimate the linkage disequilibrium between SNPs. When summary statistics were not available in full, we computed small GRS by performing clumping to select the most significant, low correlated variants, and directly using the coefficients provided in the summary statistics. Clumping and GRS computation were performed using the PLINK ([Chang et al., 2015](#)) software.

3.2.7 Statistical analyses

We estimated the association between the binary ICD phenotype and GRS using logistic regression, while correcting for age, sex, genetic ancestry (first four components), and the number of visits. We added the correction for the number of visits to reflect the fact that lifetime phenotype may be more likely as the number of visits increases. We performed the analyses in each cohort independently as the contributions of all the SNPs were estimated altogether, and the number of SNPs was much lower in PPMI than in DIGPD. We applied per-sample Bonferroni correction for multiple comparisons. We also investigated the association of the combination of the 40 GRS with the likelihood-ratio test.

As the sample sizes were relatively small in both cohorts, we also performed a meta-analysis to estimate the combined effects of each GRS separately and combined altogether using fixed effects models with the inverse-variance weighting method.

Logistic regressions were performed using the `statsmodels` (Seabold and Perktold, 2010) Python package. Meta-analyses were performed using the `meta` (Balduzzi et al., 2019) R package. Processing of the different text-like files was performed using the `pandas` (McKinney, 2010) and `NumPy` (Harris et al., 2020a) packages.

3.3 Results

3.3.1 Participants and genetic variants

Out of the 378 PPMI subjects and 382 DIGPD subjects who matched the first three inclusion criteria, we excluded 10 subjects from PPMI and 9 subjects from DIGPD for being projected too far from the European cluster (Figure 3.1). Thus, we included 368 subjects from PPMI and 373 subjects from DIGPD.

Out of the 39,235,157 genetic variants of the Haplotype Reference Consortium reference panel, 601,370 SNPs in PPMI and 6,294,320 SNPs in DIGPD matched the inclusion criteria. The high discrepancy in the numbers is due to the genotyping arrays: the NeuroX array is known to have a low coverage of the genome (Nalls et al., 2015).

3.3.2 Genome-wide association studies

We identified 40 GWAS that matched the inclusion criteria. Table 3.1 presents the characteristics of these studies, including the phenotype of interest, the number of SNPs, the heritability estimated from these SNPs, and the number of common SNPs between the GWAS and PPMI and DIGPD. The included phenotypes consisted of other psychiatric disorders (anxiety, depression, obsessive compulsive, and attention-deficit hyperactivity disorders (ADHD), anorexia nervosa), personality traits (impulsivity, agreeableness, conscientiousness, extraversion, openness), risk taking behaviors (automobile speeding, alcohol consumption, smoking status, sexual activity), and simple traits (body height, body mass index, intelligence, education).

Two groups of GWAS included genetic data from 23andMe¹, and only the top 10,000 SNPs were made publicly available. We requested access to the whole summary statistics from 23andMe with no success.

3.3.3 Association analyses

Table 3.2 presents the unadjusted p -values for the 40 GRS from the analyses on each cohort separately and from the meta-analysis. For the analysis in each cohort separately, among the 2 sets of 40 unadjusted p -values (correction is per-sample), only one was smaller than 0.05 (nominal significance), corresponding to the GRS of body mass index in PPMI ($p = 0.0079$). The association did not survive after Bonferroni correction.

¹www.23andme.com

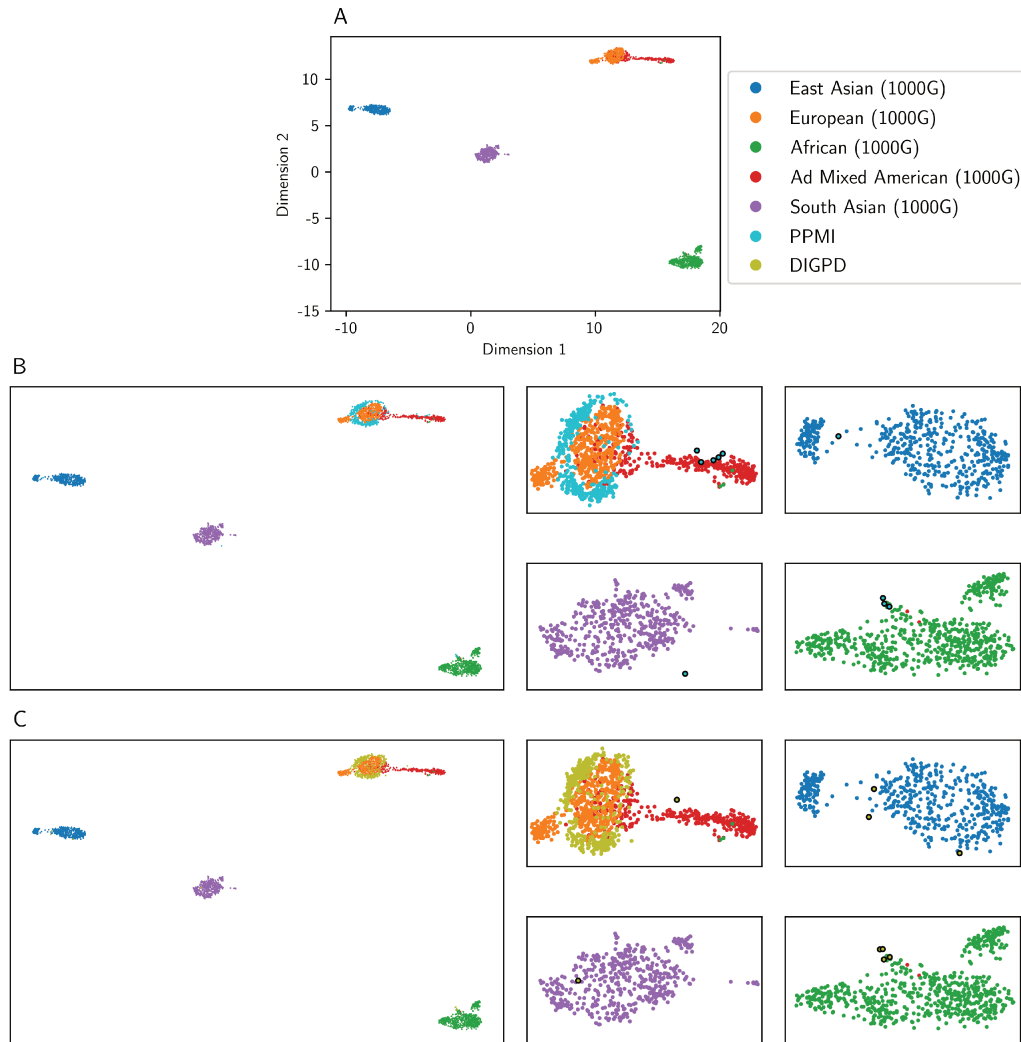


Figure 3.1: Genetic ancestry estimation. (A) A two-dimensional representation of the raw genetic data was learned on 1000G, which consists of 5 super populations. The PPMI (B) and DIGPD (C) subjects were projected on this space to estimate their genetic ancestry, and excluded if their projection was too far from the European cluster. Excluded subjects are highlighted with a black circle.

1000G: 1000 Genomes Project; DIGPD: Drug Interaction With Genes in Parkinson's Disease; PPMI: Parkinson's Progression Markers Initiative.

Study	Phenotype	h^2_{SNP} (SE)	# subjects	# SNPs	# SNPs \cap DIGPD	# SNPs \cap PPMI
(Howard et al., 2019)	Major depression disorder	0.089 (0.003)	807,553	7,743,682	5,975,580	502,668
(Nalls et al., 2019)	Parkinson’s disease	0.22 (0.024)	900,238	8,164,949	6,123,638	521,258
(IOCDF, 2018)	Obsessive compulsive disorder	0.28 (0.04)	9,725	6,813,688	5,717,746	493,510
(Otowa et al., 2016)	Anxiety disorder	0.138 (0.028)	17,31	6,330,995	5,035,054	414,627
(Watson et al., 2019)	Anorexia nervosa	0.14 (0.01)	72,517	3,448,674	2,814,932	263,385
(Karlsson Linnér et al., 2019)	Automobile speeding propensity	0.079 (0.003)	404,291	7,779,359	6,173,839	519,927
	Number of sexual partners	0.128 (0.003)	370,711	7,779,445	6,173,866	519,928
	Risk-taking tendency	0.156 (0.004)	315,894	7,779,520	6,173,864	519,925
	General risk tolerance	0.045 (0.001)	975,353	7,779,339	6,173,819	519,927
(Sanchez-Roige et al., 2019)	Drug experimentation measurement	0.1116 (0.0242)	22,572	6.442	4,850	475
	Impulsivity (attentional)	0.0541 (0.0225)	21,876	5.715	4,183	293
	Impulsivity (motor)	0.0443 (0.0193)	21,806	5.77	4,312	424
	Impulsivity (non-planning)	0.0657 (0.0254)	21,786	6.412	4,785	447
	Impulsivity	0.066 (0.0224)	21,495	5.697	4,178	268
	Lack of perseverance	0.0791 (0.019)	22,861	6.254	4,720	367
	Lack of premeditation	0.0452 (0.0199)	22,774	6.161	4,674	267
	Negative urgency	0.0796 (0.0236)	22,795	6.387	4,805	339
	Positive urgency	0.0682 (0.0233)	22,738	6.299	4,787	367
	Sensation seeking	0.0811 (0.0211)	22,745	6.769	5,342	289
(Luciano et al., 2018)	Neuroticism	0.108 (0.005)	452,688	7,625,696	6,033,335	511,477
(Savage et al., 2018)	Intelligence	0.197 (0.009)	269,867	7,445,515	6,021,486	511,723
(Demontis et al., 2019)	Attention-deficit hyperactivity disorder	0.216 (0.014)	53,293	6,921,780	5,597,529	493,284
(van den Berg et al., 2016)	Extraversion	0.050 (0.072)	72,813	6,576,855	5,411,885	452,537
(Lo et al., 2017)	Agreeableness	0.085 (0.009)	76,551	7.208	5,932	803
	Conscientiousness	0.096 (0.009)	123,132	7.267	5,777	518
	Extraversion	0.181 (0.010)	169,507	8.582	5,551	1,209
	Openness	0.107 (0.009)	76,581	7.515	5,984	447
(Pulit et al., 2019)	Body mass index	0.279 (0.002)	806,834	7,837,070	6,147,810	517,713
(Yengo et al., 2018)	Height	0.483 (0.037)	456,426	2,333,797	2,143,637	169,576
(ILAE, 2018)	Epilepsy	0.321 (0.0145)	38,752	4,988,035	4,731,775	450,99
(Liu et al., 2019)	Age of smoking initiation	0.0468 (0.0027)	341,427	7,788,606	6,139,297	520,186
	Smoking behaviour	0.0804 (0.0076)	377,334	7,788,737	6,139,286	520,188
	Smoking cessation	0.0464 (0.0018)	547,219	7,837,671	6,176,531	520,050
	Smoking initiation	0.0777 (0.0021)	1,232,091	7,683,723	6,061,508	516,602
	Alcohol consumption	0.0419 (0.0018)	941,280	7,784,169	6,136,554	519,930
(Neale lab, 2018)	Ever addicted to any substance or behaviour	0.0526 (0.0278)	26,402	8,242,335	6,017,163	511,372
	Sleeplessness / insomnia	0.0624 (0.00349)	360,738	8,247,437	6,017,509	511,412
	Trouble falling or staying asleep	0.0581 (0.00728)	117,822	8,247,396	6,017,506	511,412
	Age first had sexual intercourse	0.1614 (0.00586)	317,694	8,247,440	6,017,509	511,412
	Recent poor appetite or overeating	0.0493 (0.0074)	117,907	8,247,393	6,017,506	511,412
	Age completed full time education	0.1047 (0.00473)	240,547	8,247,414	6,017,509	511,413

Table 3.1: Characteristics of the genome-wide association studies. DIGPD: Drug Interaction With Genes in Parkinson’s Disease; PPMI: Parkinson’s Progression Markers Initiative; SE: standard error; SNP: Single nucleotide polymorphism.

	Phenotype	PPMI		DIGPD		Meta-analysis	
		OR (95% CI)	p-value	OR (95% CI)	p-value	OR (95% CI)	p-value
(Howard et al., 2019)	Major depression disorder	2.46 (0.54-11.28)	0.2449	0.64 (0.15-2.78)	0.5480	1.23 (0.43-3.53)	0.7059
(Nalls et al., 2019)	Parkinson's disease	1.47 (0.24-9.04)	0.6758	2.65 (0.77-9.12)	0.1224	2.20 (0.79-6.11)	0.1306
(IOCDF, 2018)	Obsessive compulsive disorder	4.46 (0.91-21.94)	0.0657	0.91 (0.21-3.90)	0.9039	1.88 (0.64-5.48)	0.2501
(Otowa et al., 2016)	Anxiety disorder	0.41 (0.11-1.52)	0.1829	0.74 (0.24-2.27)	0.5958	0.58 (0.25-1.35)	0.2044
(Watson et al., 2019)	Anorexia nervosa	1.20 (0.32-4.47)	0.7882	1.08 (0.32-3.60)	0.9037	1.13 (0.46-2.75)	0.7867
(Karlsson Linnér et al., 2019)	Automobile speeding propensity	2.46 (0.69-8.80)	0.1660	1.42 (0.40-5.12)	0.5885	1.87 (0.76-4.62)	0.1728
	Number of sexual partners	0.55 (0.14-2.17)	0.3914	0.54 (0.15-1.91)	0.3403	0.54 (0.21-1.38)	0.1998
	Risk-taking tendency	0.95 (0.20-4.49)	0.9482	0.99 (0.27-3.55)	0.9828	0.97 (0.36-2.61)	0.9538
	General risk tolerance	0.45 (0.10-2.07)	0.3067	1.16 (0.33-4.04)	0.8178	0.79 (0.30-2.08)	0.6368
(Sanchez-Roige et al., 2019)	Drug experimentation measurement	0.74 (0.19-2.89)	0.6691	1.32 (0.42-4.15)	0.6383	1.04 (0.43-2.49)	0.9339
	Impulsivity (attentional)	0.72 (0.16-3.16)	0.6629	2.04 (0.46-9.10)	0.3512	1.20 (0.42-3.45)	0.7297
	Impulsivity (motor)	1.04 (0.27-3.99)	0.9525	0.48 (0.14-1.58)	0.2266	0.67 (0.28-1.65)	0.3883
	Impulsivity (non-planning)	3.69 (0.83-16.43)	0.0865	1.82 (0.45-7.36)	0.4002	2.53 (0.91-7.02)	0.0742
	Impulsivity	2.47 (0.77-7.95)	0.1294	1.58 (0.37-6.71)	0.5365	2.07 (0.83-5.14)	0.1168
	Lack of perseverance	5.17 (0.96-27.86)	0.0557	0.86 (0.24-3.09)	0.8197	1.66 (0.60-4.59)	0.3295
	Lack of premeditation	0.33 (0.10-1.10)	0.0721	0.29 (0.08-1.10)	0.0689	0.31 (0.13-0.76)	0.0107
	Negative urgency	1.61 (0.44-5.89)	0.4730	1.48 (0.47-4.69)	0.5024	1.54 (0.65-3.64)	0.3281
	Positive urgency	1.33 (0.34-5.23)	0.6798	1.89 (0.59-6.06)	0.2852	1.63 (0.67-3.96)	0.2797
	Sensation seeking	0.46 (0.13-1.59)	0.2180	0.76 (0.19-3.12)	0.7062	0.57 (0.23-1.45)	0.2409
(Luciano et al., 2018)	Neuroticism	1.37 (0.42-4.47)	0.6024	1.54 (0.40-5.88)	0.5276	1.44 (0.59-3.50)	0.4188
(Savage et al., 2018)	Intelligence	0.52 (0.15-1.82)	0.3061	2.68 (0.55-13.16)	0.2248	0.97 (0.36-2.61)	0.9596
(Demontis et al., 2019)	Attention-deficit hyperactivity disorder	1.41 (0.31-6.47)	0.6578	0.32 (0.08-1.25)	0.1018	0.62 (0.22-1.71)	0.3568
(van den Berg et al., 2016)	Extraversion	0.67 (0.18-2.42)	0.5386	0.46 (0.12-1.86)	0.2796	0.57 (0.22-1.45)	0.2358
(Lo et al., 2017)	Agreeableness	2.88 (0.73-11.32)	0.1302	1.37 (0.44-4.25)	0.5857	1.85 (0.77-4.43)	0.1662
	Conscientiousness	3.12 (0.66-14.71)	0.1514	0.85 (0.23-3.20)	0.8139	1.47 (0.54-4.02)	0.4532
	Extraversion	1.06 (0.26-4.26)	0.9358	0.52 (0.16-1.71)	0.2837	0.70 (0.28-1.74)	0.4456
	Openness	1.14 (0.30-4.31)	0.8422	1.98 (0.55-7.14)	0.2981	1.52 (0.60-3.82)	0.3755
(Pulit et al., 2019)	Body mass index	23.92 (2.30-249.19)	0.0079	1.18 (0.17-8.29)	0.8696	4.04 (0.90-18.10)	0.0680
(Yengo et al., 2018)	Height	0.57 (0.06-5.09)	0.6139	1.72 (0.26-11.32)	0.5730	1.07 (0.26-4.49)	0.9215
(ILAE, 2018)	Epilepsy	0.83 (0.11-6.36)	0.8580	1.57 (0.35-6.98)	0.5566	1.25 (0.38-4.18)	0.7128
(Liu et al., 2019)	Age of smoking initiation	1.57 (0.34-7.21)	0.5590	2.03 (0.50-8.28)	0.3253	1.80 (0.64-5.07)	0.2632
	Smoking behaviour	0.48 (0.10-2.25)	0.3519	2.21 (0.43-11.20)	0.3396	0.99 (0.32-3.04)	0.9875
	Smoking cessation	1.80 (0.29-11.30)	0.5290	0.89 (0.23-3.37)	0.8626	1.13 (0.39-3.34)	0.8181
	Smoking initiation	1.31 (0.12-14.97)	0.8256	1.42 (0.35-5.84)	0.6235	1.40 (0.41-4.73)	0.5925
	Alcohol consumption	1.48 (0.29-7.43)	0.6362	0.37 (0.10-1.35)	0.1309	0.63 (0.23-1.74)	0.3770
(Neale lab, 2018)	Ever addicted to any substance or behaviour	0.88 (0.28-2.79)	0.8242	1.50 (0.35-6.33)	0.5837	1.08 (0.44-2.66)	0.8657
	Sleeplessness / insomnia	1.46 (0.35-6.13)	0.6055	1.25 (0.37-4.17)	0.7207	1.33 (0.53-3.35)	0.5444
	Trouble falling or staying asleep	1.00 (0.23-4.37)	0.9993	0.84 (0.18-3.94)	0.8255	0.92 (0.32-2.68)	0.8796
	Age first had sexual intercourse	1.11 (0.17-7.15)	0.9120	2.88 (0.83-9.99)	0.0951	2.15 (0.76-6.04)	0.1472
	Recent poor appetite or overeating	0.37 (0.10-1.39)	0.1418	2.68 (0.74-9.69)	0.1330	1.03 (0.41-2.58)	0.9572
	Age completed full time education	1.58 (0.45-5.58)	0.4783	1.00 (0.27-3.63)	0.9960	1.26 (0.51-3.11)	0.6144

Table 3.2: Results of the association analyses. Unadjusted p-values are reported for both cohorts separately and for the meta-analysis. DIGPD: Drug Interaction With Genes in Parkinson's Disease; OR: Odds ratio; PPMI: Parkinson's Progression Markers Initiative.

In the meta-analysis, among the 40 unadjusted p -values, only one was smaller than 0.05, corresponding to the GRS lack of premeditation ($p = 0.0107$). The association did not survive after Bonferroni correction.

The combination of the 40 GRS altogether was not associated with ICDs, both from the analyses in each cohort independently ($p = 0.0969$ in PPMI, $p = 0.5166$ in DIGPD) and from the meta-analysis ($p = 0.0764$).

In order to validate our GRS calculation, we assessed the quality of the GRS of body mass index by using a linear regression model with correction for age, sex and genetic ancestry. BMI was available in both cohorts, and is well studied in genetics, leading to robust GRS-based prediction (Dudbridge, 2013). In particular, the corresponding GWAS has a very large sample size ($N = 806,834$), making the estimation of each SNP contribution more robust. In both cohorts, BMI GRS were strongly associated with the measured BMI ($p = 0.000058$ in PPMI, $p = 0.000038$ in DIGPD) and the Pearson correlation coefficients were positive and high ($r = 0.21$ in PPMI, $r = 0.19$ in DIGPD). These results gave us confidence in our methodology and in the quality of the computed GRS.

3.4 Discussion

To our knowledge, this is the first study investigating the association between ICDs in PD and genetic risk scores for a broad range of phenotypes, including phenotypes that have been associated with ICDs in PD (Grall-Bronnec et al., 2018). Compared to a previous study that only investigated the PD GRS computed from a small number of SNPs (Ihle et al., 2020), we explored 40 phenotypes for which we computed GRS using a large number of SNPs. However, the results were mainly negative, as we did not find a single association after correction for multiple comparisons.

The main limitation of our study is the small sample size of our clinical samples, which limits discovery of small associations. The size of the GWAS is also a limitation, as GRS are imperfect predictors of the genetic liability of traits. It is known that discouraging results in many studies were due to low number of participants, and that an increase in the sample size would allow more successful results (Dudbridge, 2013). The genetic correlations between the traits for which GRS were calculated and ICDs are also unknown. Another limitation is the incomplete summary statistics made available for two groups of studies focusing on impulsivity and personality traits (collaborations with 23andMe, we contacted 23andMe but did not receive a response). For these traits, we had to compute GRS from a small number of SNPs. Computing the GRS using the whole summary statistics would likely increase the quality of these GRS.

Little is known about the genetic factors of ICDs in PD. Several studies reported associations for a few genetic variants, but they all suffer from the lack of replication, and there exists no GRS for ICDs yet. Our study could not conclude about the asso-

ciation between ICDs in PD and GRS for a broad range of phenotypes, but highlights the methodology to compute GRS and study their association with ICDs in PD for future studies, and shows how to investigate the genetic factors of a phenotype without performing a GWAS. Such study would deserve from being repeated when larger GWAS or clinical samples get available, which may boost power to detect significant associations.

Chapter 4

Combining static and dynamic data in recurrent neural networks

4.1 Introduction

Most entities that produce data at several time points have characteristics that do not depend on time. For instance, sensors provide measurements at many timestamps, but they have characteristics, such as their components, that are time-independent. Humans, as living beings, also have both types of characteristics. Their genetic data will stay identical for their whole lives, but most of their characteristics evolve over time, such as their blood pressure or blood sugar levels.

Numerous genetic and environmental factors can impact the evolution of a time-dependent characteristic. Twin studies, where subjects share identical genetic and environmental factors, have shown that many phenotypes are substantially heritable. Environmental factors, in particular qualitative ones, can be harder to measure, but their impact on many disorders have been reported. For example, smokers and coffee drinkers have a lower risk of Parkinson’s disease ([Hernán et al., 2002](#)).

Disease progression is a particular example of the evolution of a time-dependent characteristic. Many disorders are complex, with numerous comorbidities, and disease progression may greatly vary between patients. Being able to predict the future state of a patient may improve understanding of the disease and patient care.

Mathematically, the objective is to predict the future value of a dynamic target variable y at time $t + \tau$, $y^{(t+\tau)}$, given input features from the present and the past. These features consist of static features \mathbf{s} , that do not depend on time, and dynamic features \mathbf{x} , at time $(1, \dots, t)$, $(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)})$. Typically, the static features correspond to genetic and socio-demographic data, while the dynamic features correspond to clinical measurements. Environmental factors can be either static or dynamic feature, depending on their nature. Without loss of generality, the target variable can be included in the dynamic input features if its past and present values are relevant to predict its future

value. This machine learning task has the following mathematical formulation

$$\hat{y}^{(t+\tau)} = f\left(\mathbf{s}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}; \boldsymbol{\theta}\right)$$

where f is the function predicting $\hat{y}^{(t+\tau)}$ given the static features \mathbf{s} and input dynamic features $(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)})$, and $\boldsymbol{\theta}$ are the parameters of function f . The function f determines which information is extracted from the static and input dynamic features, independently but also dependently on each other. Different functions f allow for different modelling of the interaction between the static and input dynamic features. We restrict the choice of f to the class of recurrent neural networks (RNNs).

Recurrent neural networks are a class of neural networks dedicated to sequential data. The main part of a recurrent neural network is a RNN unit that takes as input a sequence of dynamic features and outputs a vector, corresponding to the information extracted from the dynamic features by this unit. However, a target dynamic feature is usually impacted not only by other dynamic features, but also by static features. Adding static features in recurrent neural networks raises the question of their integration with dynamic features.

4.2 Related work

The starting point is a recurrent neural network with no static data as illustrated in [Figure 4.1](#). We consider a simple architecture with three layers:

- the input layer, consisting of the dynamic features;
- the hidden layer, consisting of the features extracted by the Gated Recurrent Unit from the dynamic features; and
- the output layer, consisting of the output obtained by linearly combining the hidden layer with a Fully Connected function.

A dummy way of integrating static features is to simply remove them, and we refer to this approach as **static=none**.

Static and dynamic data can be considered as a particular combination of multi-modal data. Several studies in the medical field integrated several sources of data to improve the prediction of a phenotype. The modalities used are often imaging data, such as T1-weighted magnetic resonance imaging (T1-MRI), T2-MRI and Fluid-attenuated inversion recovery (FLAIR), and genetic data. [Ge et al. \(2018\)](#) integrated images from T1-MRI, T2-MRI and FLAIR modalities for glioma classification, while [Punjabi et al. \(2019\)](#) used T1-MRI and positron emission topography images for Alzheimer’s disease classification. [Mobadersany et al. \(2018\)](#) and [Hao et al. \(2019\)](#) integrated histopathologic images and genetic data to predict cancer outcome.

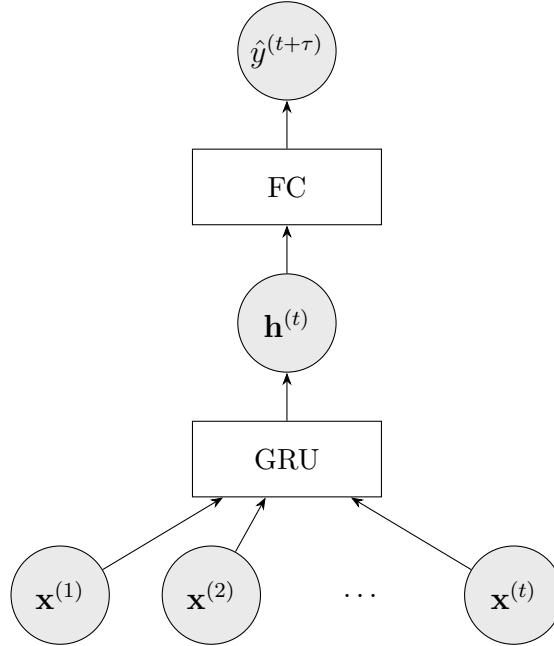


Figure 4.1: Recurrent neural network with no static data.

Integrating multimodal data is usually tackled with separate branches in the artificial neural networks, with the independently extracted features being concatenated near the end of the network. The corresponding architecture is illustrated in [Figure 4.2](#), in the particular case where the raw genetic data are used. We refer to this approach as **static=after**. One limitation of this approach is the late interaction between the static and dynamic features, as each branch extracts information from the input data independently. If the static and dynamic features are correlated, we may want to provide both kinds of features to the GRU.

Most studies focusing on the integration of static and dynamic data have been identified in the literature of churn prediction. In these studies, the objective was to predict which customers will unsubscribe to a service or which users will not log at least once into a platform in the near future. The dynamic features consisted of the activity of the users, while the static features included socio-demographic information about the users. Besides the two aforementioned methods, two other approaches have been proposed. The first one consists in treating static data as dynamic data ([Leontjeva and Kuzovkin, 2016](#); [Rahman et al., 2020](#)). The static features are repeated at each time point by being concatenated to the vector of dynamic features. This approach is illustrated in [Figure 4.3](#) and is referred to as **static=dynamic**. One obvious limitation of this method is that static data is treated as dynamic data, which may be suboptimal because of the different nature of these features.

The other approach consists in initializing the parameters of the GRU with the static features ([Kristensen and Burelli, 2019](#)). A Fully Connected function is used to

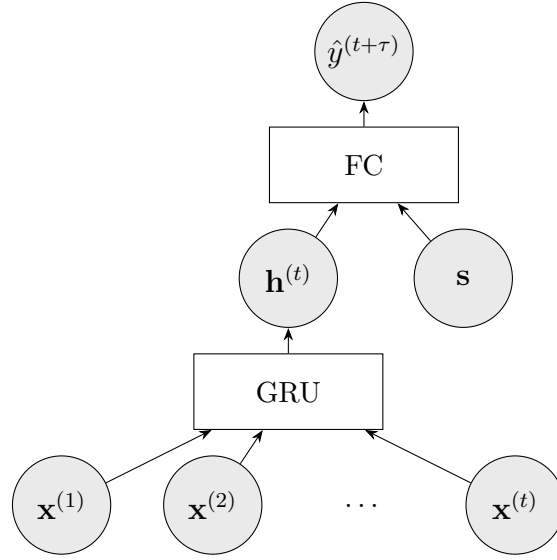


Figure 4.2: Recurrent neural network with static data on its own branch. Information is independently extracted from the dynamical and static features, in their own branches. The extracted information is then concatenated before the fully connected layer. In this example, the raw static features are directly used.

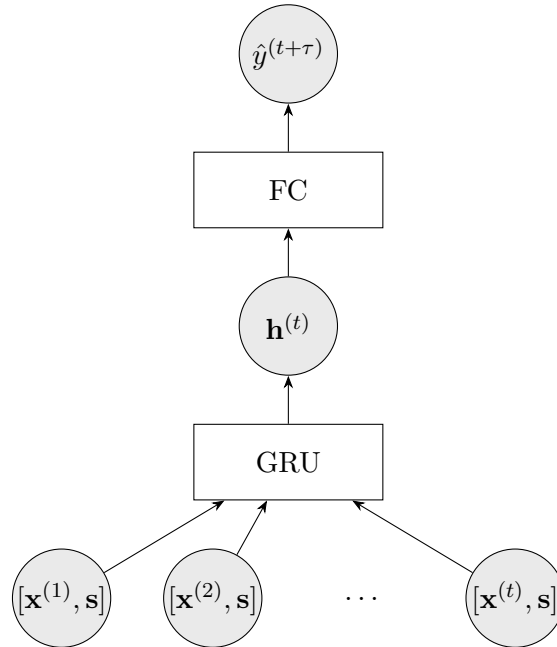


Figure 4.3: Recurrent neural network with static data treated as dynamic data. The static features are repeated at each time point and concatenated to the dynamic features.

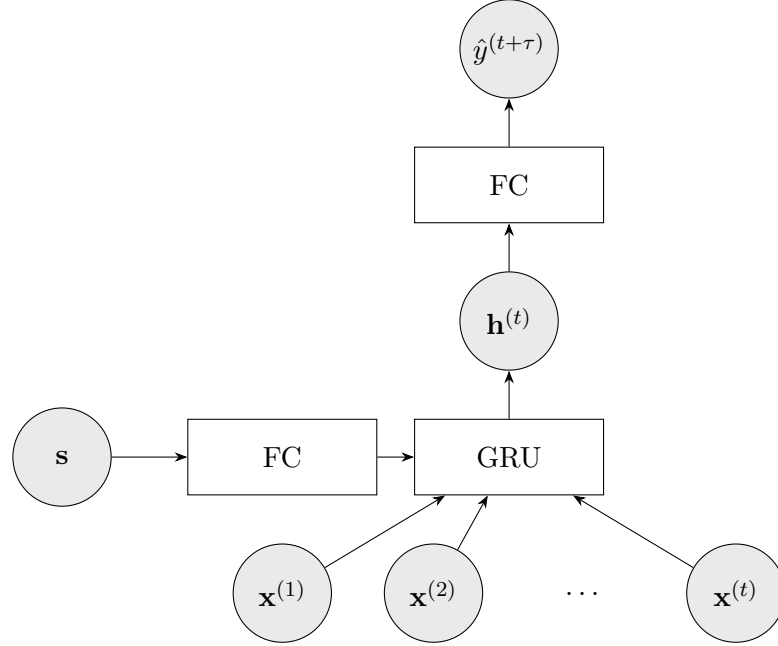


Figure 4.4: Recurrent neural network with static data initializing the GRU layer. Linear combinations of the static feature are used to initialize the parameters of the GRU layer.

linearly combine the static features into a vector with appropriate size. This approach is illustrated in [Figure 4.4](#) and is referred to as **static=init**. One limitation of this approach is that the information from the static features is added in the initialisation of the GRU and may vanish after a few time points. Otherwise, having to keep the information from the static features may prevent it from extracting other useful information from the dynamic features.

To summarize the four identified approaches, static features can be:

- removed (**static=none**),
- put after the GRU (**static=after**),
- treated as dynamic features (**static=dynamic**), or
- put at the same level as the GRU (**static=init**).

These approaches allow for modelling different interactions between the static and dynamic features.

4.3 Proposed approach

We propose another approach to integrate static data in recurrent neural networks. Similarly to the **static=dynamic** method, this approach introduces the static features

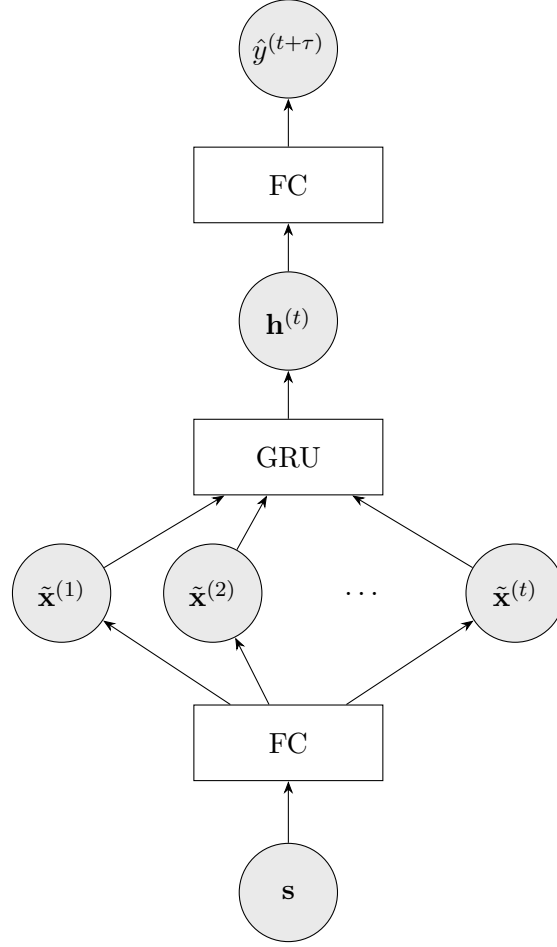


Figure 4.5: Recurrent neural network with static data modifying the dynamic features. Each dynamic feature is multiplied by a linear combination of the static features.

before the GRU. However, instead of putting the static features at the same level as the dynamic ones, they are put before and used to *modify* them.

More specifically, the new dynamic feature \tilde{x}_i is the product between the original dynamic feature x_i and a linear combination of the static features:

$$\tilde{x}_i = x_i \times \left(b + \sum_{k=1}^p s_k \right)$$

and the sequence of new dynamic features $(\tilde{x}^{(1)}, \dots, \tilde{x}^{(t)})$ is used as input of the GRU. This approach is illustrated in [Figure 4.5](#) and is referred to as **static=before**. An element-wise multiplication with a linear combination of the input boils down to adding a Fully Connected function followed by an element-wise product between the static and dynamic features.

This method models a high level of interaction between the static and dynamic features. Similarly to the **static=dynamic** and **static=init** approaches, the GRU is

provided the static features and can thus extract information from both the dynamic and static features. Instead of treating static features as dynamic ones, which does not take into account the difference in modalities, or initializing the GRU with the static features, which may dilute the information of the static features over time, this approach models a time-independent interaction between the static and dynamic features.

4.4 Experiments

We investigated the five approaches to integrate static data in recurrent neural networks to predict impulse control disorders in Parkinson’s disease. The objective was identical to the one presented in [chapter 2](#), that is predicting the presence or absence of ICDs (binary variable) at the next visit for a given subject given all the current information on this subject:

$$\hat{y}^{(t+1)} = f\left(\mathbf{s}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}; \boldsymbol{\theta}\right)$$

In addition to the clinical, socio-demographic and SNP data used, we integrated genetic ancestry and genetic risk scores as input of the models. These new features were naturally considered as static data.

We use genetic data from the 1000 Genomes project¹ to learn a low-dimensional representation of high-dimensional raw genetic data. We then projected the subjects from the PPMI and DIGPD cohorts onto this low-dimensional space and removed the subjects not being projected on the European cluster. Genetic ancestry was derived as the first ten components of this low-dimensional space and added as static features.

Genetic risk score were computed using only the common SNPs between PPMI and DIGPD. Out of the 40 GRS presented in [chapter 3](#), 13 were removed because the provided summary statistics were too small. The corresponding genome-wide association studies were performed using data from 23andMe², and only the top 10k variants were made available. [Table 4.1](#) presents the characteristics of the 27 GWAS and the corresponding GRS were added as input to the models.

The cross-validation was similar to the one presented in [chapter 2](#) and is illustrated in [Figure 4.6](#). We employed a nested cross-validation on PPMI that was used as the discovery cohort, and also evaluated the models on DIGPD that was used as the replication cohort, with 10 repetitions of the whole process. The 10 models were used to compute the mean and standard deviation for the area under the ROC curve and the average precision.

Results are presented in [Table 4.2](#). The five approaches yielded comparable results, with ROC AUC around 0.83 and 0.79 on PPMI and DIGPD respectively, and average precision around 0.53 and 0.61 on PPMI and DIGPD respectively. The `static=dynamic` method had the lowest scores on both cohorts, suggesting that this approach may be

¹<https://www.internationalgenome.org>

²<https://www.23andme.com>

Study	Phenotype	# subjects	h^2_{SNP}	# SNPs
(IOCDF, 2018)	Obsessive compulsive disorder	9,725	0.2800	495,612
(Demontis et al., 2019)	Attention-deficit hyperactivity disorder	53,293	0.2160	493,519
(Howard et al., 2019)	Major depression disorder	807,553	0.0890	499,639
(ILAE, 2018)	Epilepsy	38,752	0.3210	449,734
(Karlsson Linnér et al., 2019)	Risk-taking tendency	315,894	0.1560	516,198
	Number of sexual partners	370,711	0.1280	516,203
	Smoking status	518,633	0.1090	516,204
	General risk tolerance	975,353	0.0450	516,205
	Automobile speeding propensity	404,291	0.0790	516,206
	Alcohol consumption	414,343	0.0850	516,200
(Liu et al., 2019)	Smoking behaviour	377,334	0.0800	516,112
	Alcohol consumption	941,280	0.0420	515,914
	Smoking initiation	1,232,091	0.0780	512,648
	Age of smoking initiation	341,427	0.0470	516,112
	Smoking cessation	547,219	0.0460	515,985
(Luciano et al., 2018)	Neuroticism	452,688	0.1080	518,481
(Nalls et al., 2019)	Parkinson's disease	900,238	0.2600	518,569
(Neale lab, 2018)	Trouble falling or staying asleep	117,822	0.0581	517,205
	Sleeplessness / insomnia	360,738	0.0624	517,204
	Age first had sexual intercourse	317,694	0.1614	517,205
	Ever addicted to any substance or behaviour	26,402	0.0526	516,916
	Age completed full time education	240,547	0.1047	517,205
	Recent poor appetite or overeating	117,907	0.0493	517,205
(Otowa et al., 2016)	Anxiety disorder	17,31	0.1380	413,915
(Pulit et al., 2019)	Body mass index	806,834	0.2790	513,153
(Savage et al., 2018)	Intelligence	269,867	0.2050	507,254
(Watson et al., 2019)	Anorexia nervosa	72,517	0.1400	263,451
(Yengo et al., 2018)	Height	456,426	0.4830	169,402
(van den Berg et al., 2016)	Extraversion	72,813	0.0500	452,260

Table 4.1: Genome-wide association studies from which genetic risk scores were derived. Columns are: study, phenotype, number of subjects in the study, variance explained by the SNPs, and the number of common SNPs between the study, ADNI, PPMI and DIGPD.

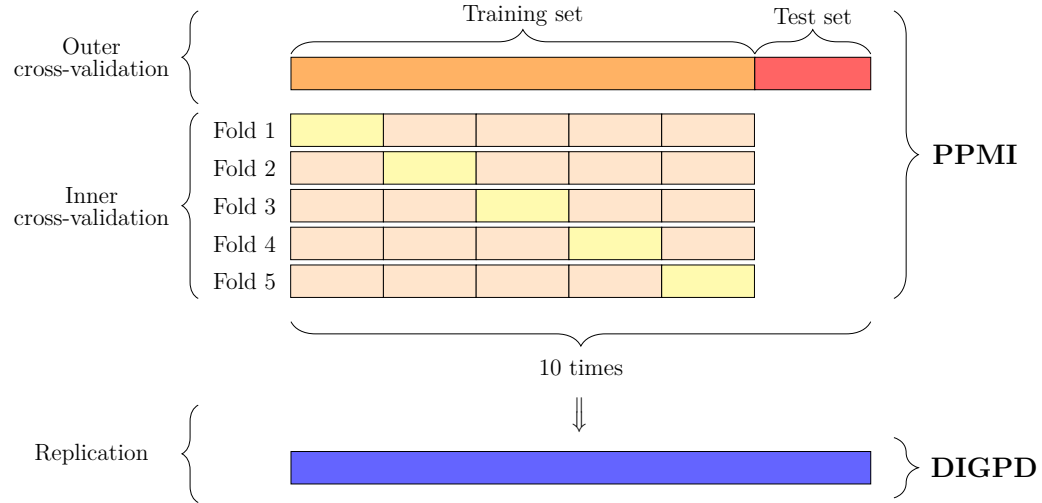


Figure 4.6: Cross-validation procedure. PPMI was used as the discovery cohort, on which we performed a nested cross-validation to estimate the performance. This process was repeated 10 times, and the corresponding 10 models were applied on DIGPD, used as the replication cohort.

Method	ROC AUC		Average Precision	
	PPMI	DIGPD	PPMI	DIGPD
<code>static=none</code>	0.837 (± 0.027)	0.788 (± 0.007)	0.543 (± 0.041)	0.614 (± 0.013)
<code>static=after</code>	0.832 (± 0.031)	0.789 (± 0.007)	0.520 (± 0.047)	0.612 (± 0.014)
<code>static=dynamic</code>	0.816 (± 0.038)	0.782 (± 0.008)	0.494 (± 0.063)	0.606 (± 0.015)
<code>static=init</code>	0.837 (± 0.033)	0.793 (± 0.007)	0.537 (± 0.047)	0.619 (± 0.011)
<code>static=before</code>	0.839 (± 0.030)	0.800 (± 0.007)	0.536 (± 0.054)	0.635 (± 0.016)

Table 4.2: Predictive performance of the five approaches. Mean (standard deviation) over the 10 repetitions are reported.

suboptimal. Overall, the different approaches to integrate static data in recurrent neural networks had little to no impact on the predictive performance with these data.

4.5 Conclusion

Combining static and dynamic data is an underexplored topic. We reviewed the literature on integrating static data in recurrent neural networks and identified four approaches. We proposed a new method modelling a high-level interaction between the static and dynamic features, consisting in multiplying each dynamic feature by a linear combination of the static features. We experimented the five approaches in the prediction of future impulse control disorders in Parkinson’s disease. The dynamic features consisted in clinical measurements, while the static features corresponded to genetic

and socio-demographics data. The results showed that the five approaches yielded comparable results. This use case has been underexplored from a machine learning point of view and little is known on the association between genetics and impulse control disorders in Parkinson’s disease, which could explain why changing the approach had little to no impact on the predictive performance. Future work includes applying the five approaches to another use case, where the interaction between the static and dynamic features is known to be high, and simulating data to gather more knowledge on which approach works best based on the interaction.

Conclusion and perspectives

Conclusion

We proposed several approaches to investigate the predictability of impulse control disorders in Parkinson’s disease. First, we investigated the predictability of ICDs in PD using as input a broad range of features that have been associated with ICDs. Second, we studied the association between ICDs in PD and genetic risk scores for numerous phenotypes, including other psychiatric disorders and personality traits. Third, we investigated the integration of static (time-independent) features in recurrent neural networks with an application in the prediction of ICDs in PD. We summarize here our conclusions regarding each of these studies.

In a first study, we investigated the predictability of ICDs in PD using as input a broad range of features that have been associated with ICDs (with varying degrees of confidence and replication). Our objective was to predict the presence or absence of ICDs at the next visit for a given patient. We trained several machine learning algorithms, representing a broad range of complexity and relationships, on a discovery cohort. We highlighted that this longitudinal binary classification task could be addressed with relatively good accuracy, and that only a subset of the associated factors was involved in the decisions of the simplest models. We also evaluated the models in an independent cohort and obtained comparable results.

In a second study, we investigated whether the genetic factors of many traits were associated with ICDs in PD. The genetic factors of ICDs in PD are poorly known and most studies focus on candidate genetic variants. However, complex traits are likely to be affected by numerous variants and genes. In this case, single genetic variants usually provide limited information because the relative risk of each variant is small. On the other hand, the combined risk of numerous low-risk variants can explain a significant proportion of the genetic variance. The risk of each variant is linearly combined to derive a genetic risk score. We computed genetic risk scores for a broad range of phenotypes, including other psychiatric disorders, personality traits, and simple phenotypes. We assessed the associations between these genetic risk scores and ICDs in PD and found

no association.

In a third study, we investigated the integration of static features in recurrent neural networks. We reviewed the existing literature and identified several approaches. We proposed a new approach consisting in modifying the dynamic features using a linear combination (or more generally a function) of the static features, modelling a high level of interaction between the static and dynamic features. We evaluated all the approaches in the use case of predicting the presence or absence of ICDs at the next visit for a given patient. The static features consisted of socio-demographic and genetic features. All the approaches (including removing the static features) led to very similar results, suggesting that the static features do not provide more information than the dynamic features in this use case.

Perspectives

A natural perspective when developing machine learning models for a medical application is the increase of sample size.

First, a larger sample size on the training cohort can improve the generalization capability of the fitted models. Learning curves in machine learning usually show that, when the sample size is relatively small, an increase in sample size leads to a significant improvement in predictive performance, but when the sample size is already large enough, a similar increase in sample size has little to no impact on the predictive performance. Learning curves, and thus the definition of low and large sample sizes, highly depend on the algorithm and the difficulty of the task.

Second, a larger sample size increases the statistical power and can allow for the discovery of associations with small effects. Both research cohorts from which we obtained data had only a few hundreds subjects, making the discovery of associations with small effects very unlikely. A larger sample size could shed a new light on the common genetic factors between ICDs in PD and other phenotypes, notably other psychiatric disorders and personality traits.

Third, we could evaluate our models in other research cohorts. We showed that the replication on an independent cohort with different characteristics was possible. Evaluating our models in other research cohorts would be a big step towards estimating their generalization capability.

Fourth, we did not evaluate our models in real-life clinical cohorts, which is an essential step to assess the generalization capability of the models and to deploy the models

in clinical routine. Clinical cohorts usually have different characteristics from research cohorts: more missing data, less clean data, more varying time gaps between consecutive visits, less standardized protocols. All these differences can negatively impact the predictive performance of models that have not been trained on data with these characteristics.

Fifth, we showed that predicting impulse control disorders at the next clinical visit can be achieved with correct accuracy, but the predictive performance of these models may still be too low to be used in clinical routine. Investigating which threshold to binarize the predicted probabilities into decisions is the most adapted in practice would be of great interest. Combining the predictions of the models and the expertise of clinicians may also improve the predictive performance.

Finally, impulse control disorders have been reported in a few other diseases treated with dopamine replacement therapy, notably restless leg syndrome. Evaluating our models in cohorts with another disease could be interesting in order to see if ICDs in PD have different characteristics than ICDs in other diseases.

Appendix A

Supplementary materials for the prediction of impulse control disorders from clinical and genetic data with replication in an independent cohort

A.1 Reduction approaches

Algorithms like logistic regression expect a fixed number of features as input. In order to deal with varying numbers of visits, we reduced all the previous visits into one “summary” visit using a convex combination. A convex combination is a linear combination such that the weights are all non-negative and sum to one. The weights indicate how much each visit contributes to this summary visit. A weight of 1 for the first visit means that the summary visit is simply the baseline visit, while a weight of 1 for the latest visit means that the “summary” visit is simply the most recent visit. One can also give uniform weights, so that each visit contributes equally to this “summary” visit, or higher weights to most recent visits if they are assumed to be more important than older visits.

Mathematically, if we have observations \mathbf{x} at T time points:

$$(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(T)})$$

a convex combination is simply:

$$\sum_{t=1}^T w_t \mathbf{x}^{(t)} \quad \text{such that} \quad \forall t, w_t \geq 0 \quad \text{and} \quad \sum_{t=1}^T w_t = 1$$

Each w_t is the weight of time point t in this “summary” visit. The following table summarizes the different convex combinations that we investigated:

Name	Weight w_t
Reduction 1	$\forall t \in \{1, \dots, T\}, w_t = \begin{cases} 1 & \text{if } t = 1 \\ 0 & \text{otherwise} \end{cases}$
Reduction 2	$\forall t \in \{1, \dots, T\}, w_t = \begin{cases} 1 & \text{if } t = T \\ 0 & \text{otherwise} \end{cases}$
Reduction 3	$\forall t \in \{1, \dots, T\}, w_t = \frac{1}{T}$
Reduction 4	$\forall t \in \{1, \dots, T\}, w_t = \frac{\exp(\sqrt{t})}{\sum_{l=1}^T \exp(\sqrt{l})}$
Reduction 5	$\forall t \in \{1, \dots, T\}, w_t = \frac{\sqrt{t}}{\sum_{l=1}^T \sqrt{l}}$
Reduction 6	$\forall t \in \{1, \dots, T\}, w_t = \frac{\exp(t)}{\sum_{l=1}^T \exp(l)}$
Reduction 7	$\forall t \in \{1, \dots, T\}, w_t = \frac{t}{\sum_{l=1}^T l}$
Reduction 8	$\forall t \in \{1, \dots, T\}, w_t = \frac{\exp(t^2)}{\sum_{l=1}^T \exp(l^2)}$
Reduction 9	$\forall t \in \{1, \dots, T\}, w_t = \frac{t^2}{\sum_{l=1}^T l^2}$

Reduction 1 corresponds to the baseline visit, while reduction 2 corresponds to the previous visit, and reduction 3 corresponds to the mean over the past visits.

A.2 Supplementary Tables

Gene	rs
<i>ANKK1</i>	rs1800497
<i>ARC</i>	rs10097505
<i>BDNF</i>	rs6265
<i>C8B</i>	rs617283
<i>C8B</i>	rs725330
<i>C8B</i>	rs591730
<i>CA12</i>	rs1075456
<i>CA12</i>	rs7166946
<i>CA12</i>	rs1043239
<i>CA12</i>	rs4984241
<i>CA12</i>	rs1043256
<i>CA12</i>	rs9989288
<i>CA12</i>	rs2046484
<i>CA12</i>	rs16946963
<i>CCRN4L</i>	rs938836
<i>COMT</i>	rs4680
<i>DBH</i>	rs1108580
<i>DBH</i>	rs1611115
<i>DRD2</i>	rs6277
<i>DRD3</i>	rs6280
<i>FOSB</i>	rs1049739
<i>FOSB</i>	rs2282695
<i>FOSB</i>	rs2276469
<i>GRIN2B</i>	rs1806201
<i>HTR1B</i>	rs6296
<i>MOSC1</i>	rs1109103
<i>MOSC1</i>	rs2984657
<i>OPRM1</i>	rs1799971
<i>TPH1</i>	rs1800532
<i>TPH2</i>	rs1352250
<i>DBH</i>	rs1108580

Table A.1: Genetic variants included in the analyses.

		Baseline visit	Most recent visit	Mean over past visits
Socio-demographic	Sex	0.205	0.073	0.168
	Age	-0.309	-0.299	-0.221
Clinical	Past ICDs	1.735	2.125	3.617
	Depression	0.000	0.000	0.182
	State anxiety	0.000	0.000	0.000
	Trait anxiety	0.000	0.000	0.000
	REM sleep	0.721	0.461	0.546
	Motor exam	0.000	0.000	0.000
PD medication	On levodopa	0.000	0.000	0.000
	On dopamine agonists	0.000	0.000	0.200
	On other PD medication	0.000	0.069	0.000
	Mean daily dose of dopamine agonists	0.000	0.000	0.000
	Maximum daily dose of dopamine agonists	0.000	0.074	0.094
	Total dose of dopamine agonists	0.000	0.000	0.000
	Cumulative duration on dopamine agonists	0.000	0.116	0.000
Time to prediction	Time to prediction	0.117	0.036	0.052
Genetic	rs10097505	-0.240	0.000	-0.267
	rs1043239	0.000	0.000	0.000
	rs1043256	0.000	0.030	0.000
	rs1049739	0.000	0.000	0.000
	rs1075456	0.000	0.000	0.000
	rs1108580	0.000	0.000	0.000
	rs1109103	0.000	0.000	0.000
	rs1352250	0.000	0.000	0.000
	rs1611115	0.000	0.000	0.000
	rs16946963	0.000	0.000	-0.047
	rs1799971	-0.106	-0.051	0.000
	rs1800497	0.000	0.000	0.001
	rs1800532	-0.086	0.000	-0.075
	rs1806201	0.050	0.000	0.037
	rs2046484	0.000	0.000	0.000
	rs2276469	-0.232	-0.096	0.000
	rs2282695	0.427	0.000	0.000
	rs2984657	0.000	0.000	0.000
	rs4680	0.000	0.000	0.000
	rs4984241	0.000	0.000	0.000
	rs591730	0.000	0.000	0.000
	rs617283	0.000	0.000	-0.082
	rs6265	0.030	0.000	0.000
	rs6277	-0.109	0.000	-0.100
	rs6280	0.000	0.000	-0.033
	rs6296	-0.044	-0.018	0.000
	rs6582078	-0.120	-0.009	-0.024
	rs7166946	0.000	-0.044	-0.093
	rs725330	-0.210	-0.090	-0.023
	rs938836	0.365	0.164	0.234
	rs9989288	0.000	0.000	0.000

Table A.2: Coefficients of the three logistic regression models with genetic variants as input.

Metric	Algorithm	Reduction								
		1	2	3	4	5	6	7	8	9
ROC AUC	LogisticRegression	0.753	0.795	0.838	0.839	0.793	0.832	0.794	0.809	0.794
	LinearSVC	0.754	0.786	0.827	0.828	0.815	0.823	0.816	0.795	0.815
	SVC	0.757	0.809	0.840	0.841	0.835	0.838	0.836	0.826	0.836
	RandomForestClassifier	0.652	0.765	0.797	0.809	0.803	0.789	0.808	0.811	0.797
	XGBClassifier	0.696	0.825	0.822	0.842	0.829	0.841	0.836	0.834	0.839
Average precision	LogisticRegression	0.441	0.453	0.603	0.589	0.491	0.558	0.495	0.461	0.495
	LinearSVC	0.425	0.457	0.588	0.582	0.531	0.553	0.535	0.468	0.530
	SVC	0.509	0.504	0.628	0.603	0.570	0.589	0.573	0.521	0.574
	RandomForestClassifier	0.367	0.397	0.501	0.488	0.503	0.459	0.509	0.484	0.471
	XGBClassifier	0.399	0.448	0.476	0.502	0.506	0.459	0.547	0.476	0.539
Accuracy	LogisticRegression	0.774	0.819	0.841	0.847	0.819	0.843	0.824	0.837	0.819
	LinearSVC	0.761	0.841	0.854	0.860	0.862	0.849	0.856	0.841	0.845
	SVC	0.813	0.877	0.875	0.869	0.873	0.877	0.873	0.877	0.877
	RandomForestClassifier	0.873	0.776	0.849	0.856	0.854	0.849	0.869	0.847	0.862
	XGBClassifier	0.609	0.697	0.798	0.815	0.809	0.798	0.787	0.794	0.789
Balanced accuracy	LogisticRegression	0.691	0.764	0.765	0.775	0.770	0.761	0.773	0.763	0.758
	LinearSVC	0.683	0.771	0.767	0.777	0.784	0.753	0.780	0.748	0.756
	SVC	0.696	0.752	0.791	0.782	0.778	0.757	0.773	0.752	0.775
	RandomForestClassifier	0.584	0.639	0.659	0.668	0.667	0.665	0.676	0.663	0.678
	XGBClassifier	0.629	0.751	0.775	0.779	0.752	0.763	0.745	0.767	0.758
Sensitivity	LogisticRegression	0.571	0.686	0.657	0.671	0.700	0.643	0.700	0.657	0.671
	LinearSVC	0.571	0.671	0.643	0.657	0.671	0.614	0.671	0.614	0.629
	SVC	0.529	0.571	0.671	0.657	0.643	0.586	0.629	0.571	0.629
	RandomForestClassifier	0.171	0.443	0.386	0.400	0.400	0.400	0.400	0.400	0.414
	XGBClassifier	0.657	0.829	0.743	0.729	0.671	0.714	0.686	0.729	0.714
Specificity	LogisticRegression	0.810	0.843	0.873	0.878	0.841	0.878	0.846	0.868	0.846
	LinearSVC	0.795	0.871	0.891	0.896	0.896	0.891	0.889	0.881	0.884
	SVC	0.863	0.932	0.911	0.906	0.914	0.929	0.916	0.932	0.922
	RandomForestClassifier	0.997	0.835	0.932	0.937	0.934	0.929	0.952	0.927	0.942
	XGBClassifier	0.600	0.673	0.808	0.830	0.833	0.813	0.805	0.805	0.803

Table A.3: Predictive performance on DIGPD of the five machine learning algorithms with the nine reduction approaches.

Metric	Algorithm	Reduction								
		1	2	3	4	5	6	7	8	9
ROC AUC	LogisticRegression	0.666	0.802	0.797	0.811	0.817	0.813	0.819	0.804	0.821
	LinearSVC	0.679	0.802	0.796	0.812	0.809	0.815	0.811	0.805	0.812
	SVC	0.623	0.784	0.758	0.786	0.780	0.800	0.787	0.791	0.790
	RandomForestClassifier	0.593	0.764	0.751	0.737	0.717	0.734	0.715	0.728	0.715
	XGBClassifier	0.639	0.791	0.774	0.796	0.789	0.799	0.798	0.796	0.802
Average precision	LogisticRegression	0.429	0.644	0.615	0.643	0.635	0.650	0.636	0.644	0.640
	LinearSVC	0.420	0.641	0.611	0.636	0.639	0.642	0.644	0.640	0.645
	SVC	0.389	0.614	0.579	0.624	0.618	0.634	0.633	0.623	0.626
	RandomForestClassifier	0.339	0.605	0.512	0.537	0.509	0.523	0.528	0.527	0.541
	XGBClassifier	0.394	0.652	0.540	0.618	0.609	0.629	0.639	0.639	0.649
Accuracy	LogisticRegression	0.566	0.592	0.645	0.666	0.500	0.643	0.511	0.609	0.515
	LinearSVC	0.590	0.582	0.658	0.688	0.652	0.674	0.651	0.624	0.647
	SVC	0.566	0.838	0.696	0.759	0.719	0.797	0.760	0.838	0.748
	RandomForestClassifier	0.719	0.711	0.757	0.764	0.753	0.746	0.751	0.763	0.762
	XGBClassifier	0.546	0.630	0.738	0.769	0.748	0.764	0.761	0.760	0.756
Balanced accuracy	LogisticRegression	0.603	0.668	0.695	0.708	0.634	0.696	0.640	0.676	0.641
	LinearSVC	0.621	0.661	0.705	0.728	0.703	0.717	0.702	0.689	0.700
	SVC	0.588	0.783	0.703	0.752	0.727	0.768	0.754	0.783	0.746
	RandomForestClassifier	0.501	0.704	0.632	0.608	0.614	0.562	0.584	0.610	0.604
	XGBClassifier	0.588	0.683	0.746	0.764	0.746	0.759	0.751	0.759	0.750
Sensitivity	LogisticRegression	0.684	0.838	0.805	0.801	0.929	0.814	0.924	0.823	0.920
	LinearSVC	0.690	0.838	0.810	0.816	0.814	0.812	0.814	0.831	0.818
	SVC	0.636	0.662	0.719	0.736	0.745	0.703	0.740	0.662	0.742
	RandomForestClassifier	0.022	0.688	0.357	0.262	0.307	0.154	0.216	0.273	0.253
	XGBClassifier	0.682	0.801	0.764	0.753	0.742	0.749	0.729	0.758	0.738
Specificity	LogisticRegression	0.522	0.499	0.584	0.616	0.338	0.578	0.356	0.529	0.362
	LinearSVC	0.552	0.485	0.601	0.639	0.591	0.622	0.590	0.546	0.582
	SVC	0.540	0.905	0.688	0.768	0.710	0.833	0.768	0.905	0.750
	RandomForestClassifier	0.981	0.719	0.907	0.954	0.921	0.970	0.952	0.948	0.954
	XGBClassifier	0.495	0.565	0.728	0.775	0.750	0.769	0.773	0.761	0.763

Table A.4: Predictive performance on DIGPD of the five machine learning algorithms with the nine reduction approaches.

Metric	Algorithm	Reduction								
		1	2	3	4	5	6	7	8	9
ROC AUC	LogisticRegression	0.757 (\pm 0.038)	0.773 (\pm 0.035)	0.821 (\pm 0.020)	0.825 (\pm 0.019)	0.807 (\pm 0.025)	0.816 (\pm 0.025)	0.812 (\pm 0.024)	0.783 (\pm 0.032)	0.807 (\pm 0.025)
	LinearSVC	0.758 (\pm 0.040)	0.779 (\pm 0.035)	0.821 (\pm 0.016)	0.824 (\pm 0.019)	0.811 (\pm 0.022)	0.825 (\pm 0.019)	0.813 (\pm 0.024)	0.792 (\pm 0.033)	0.815 (\pm 0.022)
	SVC	0.696 (\pm 0.051)	0.767 (\pm 0.030)	0.826 (\pm 0.022)	0.828 (\pm 0.022)	0.816 (\pm 0.022)	0.816 (\pm 0.022)	0.811 (\pm 0.030)	0.788 (\pm 0.030)	0.805 (\pm 0.030)
	RandomForestClassifier	0.664 (\pm 0.033)	0.740 (\pm 0.036)	0.800 (\pm 0.028)	0.802 (\pm 0.026)	0.795 (\pm 0.032)	0.803 (\pm 0.032)	0.795 (\pm 0.030)	0.798 (\pm 0.032)	0.791 (\pm 0.030)
	XGBClassifier	0.691 (\pm 0.044)	0.782 (\pm 0.031)	0.815 (\pm 0.017)	0.832 (\pm 0.020)	0.820 (\pm 0.029)	0.829 (\pm 0.022)	0.821 (\pm 0.027)	0.824 (\pm 0.017)	0.823 (\pm 0.023)
Average precision	LogisticRegression	0.382 (\pm 0.083)	0.458 (\pm 0.064)	0.503 (\pm 0.077)	0.515 (\pm 0.064)	0.490 (\pm 0.059)	0.514 (\pm 0.062)	0.499 (\pm 0.057)	0.471 (\pm 0.067)	0.499 (\pm 0.054)
	LinearSVC	0.386 (\pm 0.084)	0.468 (\pm 0.051)	0.506 (\pm 0.075)	0.519 (\pm 0.063)	0.502 (\pm 0.062)	0.520 (\pm 0.065)	0.504 (\pm 0.060)	0.479 (\pm 0.051)	0.506 (\pm 0.057)
	SVC	0.336 (\pm 0.094)	0.458 (\pm 0.061)	0.511 (\pm 0.072)	0.528 (\pm 0.065)	0.514 (\pm 0.064)	0.532 (\pm 0.061)	0.514 (\pm 0.065)	0.486 (\pm 0.063)	0.515 (\pm 0.060)
	RandomForestClassifier	0.294 (\pm 0.059)	0.403 (\pm 0.092)	0.473 (\pm 0.067)	0.495 (\pm 0.065)	0.494 (\pm 0.066)	0.466 (\pm 0.055)	0.492 (\pm 0.070)	0.454 (\pm 0.071)	0.465 (\pm 0.056)
	XGBClassifier	0.338 (\pm 0.086)	0.455 (\pm 0.057)	0.477 (\pm 0.056)	0.529 (\pm 0.045)	0.513 (\pm 0.056)	0.507 (\pm 0.051)	0.527 (\pm 0.050)	0.479 (\pm 0.059)	0.510 (\pm 0.046)
Accuracy	LogisticRegression	0.758 (\pm 0.026)	0.836 (\pm 0.036)	0.832 (\pm 0.030)	0.846 (\pm 0.032)	0.840 (\pm 0.028)	0.849 (\pm 0.025)	0.844 (\pm 0.030)	0.845 (\pm 0.029)	0.845 (\pm 0.030)
	LinearSVC	0.765 (\pm 0.031)	0.849 (\pm 0.033)	0.847 (\pm 0.025)	0.856 (\pm 0.025)	0.856 (\pm 0.028)	0.857 (\pm 0.027)	0.855 (\pm 0.028)	0.855 (\pm 0.031)	0.854 (\pm 0.030)
	SVC	0.762 (\pm 0.042)	0.872 (\pm 0.022)	0.850 (\pm 0.027)	0.858 (\pm 0.026)	0.858 (\pm 0.031)	0.871 (\pm 0.021)	0.863 (\pm 0.029)	0.873 (\pm 0.021)	0.863 (\pm 0.027)
	RandomForestClassifier	0.827 (\pm 0.043)	0.811 (\pm 0.034)	0.849 (\pm 0.041)	0.856 (\pm 0.026)	0.857 (\pm 0.037)	0.850 (\pm 0.029)	0.852 (\pm 0.040)	0.846 (\pm 0.027)	0.846 (\pm 0.038)
	XGBClassifier	0.610 (\pm 0.067)	0.691 (\pm 0.070)	0.788 (\pm 0.023)	0.797 (\pm 0.013)	0.784 (\pm 0.028)	0.789 (\pm 0.019)	0.776 (\pm 0.026)	0.789 (\pm 0.019)	0.778 (\pm 0.015)
Balanced accuracy	LogisticRegression	0.690 (\pm 0.043)	0.734 (\pm 0.040)	0.773 (\pm 0.017)	0.782 (\pm 0.030)	0.774 (\pm 0.026)	0.770 (\pm 0.033)	0.777 (\pm 0.030)	0.748 (\pm 0.030)	0.761 (\pm 0.030)
	LinearSVC	0.691 (\pm 0.045)	0.735 (\pm 0.036)	0.774 (\pm 0.016)	0.780 (\pm 0.024)	0.775 (\pm 0.028)	0.762 (\pm 0.025)	0.774 (\pm 0.028)	0.744 (\pm 0.033)	0.760 (\pm 0.026)
	SVC	0.633 (\pm 0.046)	0.738 (\pm 0.026)	0.768 (\pm 0.022)	0.771 (\pm 0.024)	0.761 (\pm 0.018)	0.751 (\pm 0.022)	0.758 (\pm 0.024)	0.739 (\pm 0.027)	0.745 (\pm 0.026)
	RandomForestClassifier	0.581 (\pm 0.035)	0.649 (\pm 0.048)	0.684 (\pm 0.038)	0.690 (\pm 0.027)	0.687 (\pm 0.032)	0.692 (\pm 0.033)	0.674 (\pm 0.037)	0.687 (\pm 0.037)	0.675 (\pm 0.039)
	XGBClassifier	0.627 (\pm 0.045)	0.706 (\pm 0.036)	0.774 (\pm 0.016)	0.773 (\pm 0.017)	0.753 (\pm 0.032)	0.767 (\pm 0.016)	0.754 (\pm 0.033)	0.774 (\pm 0.015)	0.757 (\pm 0.021)
Sensitivity	LogisticRegression	0.593 (\pm 0.101)	0.591 (\pm 0.065)	0.689 (\pm 0.031)	0.693 (\pm 0.053)	0.683 (\pm 0.050)	0.659 (\pm 0.073)	0.684 (\pm 0.057)	0.613 (\pm 0.057)	0.644 (\pm 0.064)
	LinearSVC	0.587 (\pm 0.106)	0.574 (\pm 0.057)	0.673 (\pm 0.039)	0.673 (\pm 0.047)	0.662 (\pm 0.050)	0.630 (\pm 0.050)	0.661 (\pm 0.051)	0.588 (\pm 0.055)	0.630 (\pm 0.052)
	SVC	0.451 (\pm 0.106)	0.549 (\pm 0.051)	0.654 (\pm 0.046)	0.648 (\pm 0.047)	0.624 (\pm 0.035)	0.584 (\pm 0.047)	0.609 (\pm 0.048)	0.551 (\pm 0.053)	0.580 (\pm 0.053)
	RandomForestClassifier	0.231 (\pm 0.095)	0.419 (\pm 0.113)	0.449 (\pm 0.089)	0.453 (\pm 0.073)	0.447 (\pm 0.080)	0.465 (\pm 0.073)	0.422 (\pm 0.097)	0.460 (\pm 0.087)	0.432 (\pm 0.098)
	XGBClassifier	0.656 (\pm 0.124)	0.727 (\pm 0.072)	0.753 (\pm 0.028)	0.739 (\pm 0.034)	0.711 (\pm 0.051)	0.737 (\pm 0.034)	0.722 (\pm 0.048)	0.754 (\pm 0.031)	0.728 (\pm 0.033)
Specificity	LogisticRegression	0.786 (\pm 0.035)	0.877 (\pm 0.036)	0.856 (\pm 0.033)	0.871 (\pm 0.033)	0.866 (\pm 0.029)	0.880 (\pm 0.029)	0.870 (\pm 0.030)	0.883 (\pm 0.030)	0.878 (\pm 0.033)
	LinearSVC	0.796 (\pm 0.039)	0.896 (\pm 0.033)	0.876 (\pm 0.027)	0.886 (\pm 0.025)	0.888 (\pm 0.026)	0.895 (\pm 0.026)	0.887 (\pm 0.026)	0.900 (\pm 0.030)	0.891 (\pm 0.029)
	SVC	0.816 (\pm 0.050)	0.926 (\pm 0.016)	0.883 (\pm 0.029)	0.894 (\pm 0.023)	0.898 (\pm 0.028)	0.919 (\pm 0.018)	0.906 (\pm 0.024)	0.927 (\pm 0.016)	0.911 (\pm 0.022)
	RandomForestClassifier	0.930 (\pm 0.046)	0.879 (\pm 0.039)	0.919 (\pm 0.041)	0.928 (\pm 0.030)	0.928 (\pm 0.036)	0.919 (\pm 0.028)	0.927 (\pm 0.048)	0.914 (\pm 0.027)	0.918 (\pm 0.045)
	XGBClassifier	0.597 (\pm 0.087)	0.685 (\pm 0.086)	0.794 (\pm 0.026)	0.806 (\pm 0.017)	0.796 (\pm 0.028)	0.798 (\pm 0.023)	0.785 (\pm 0.026)	0.795 (\pm 0.023)	0.786 (\pm 0.015)

Table A.5: Predictive performance on DIGPD of the five machine learning algorithms with the nine reduction approaches with 10 repetitions of the nest cross-validation. Mean (standard deviation) over the 10 models are reported.

Metric	Algorithm	Reduction								
		1	2	3	4	5	6	7	8	9
ROC AUC	LogisticRegression	0.653 (\pm 0.033)	0.796 (\pm 0.026)	0.780 (\pm 0.031)	0.804 (\pm 0.021)	0.796 (\pm 0.032)	0.811 (\pm 0.019)	0.804 (\pm 0.024)	0.801 (\pm 0.023)	0.806 (\pm 0.022)
	LinearSVC	0.666 (\pm 0.040)	0.792 (\pm 0.023)	0.784 (\pm 0.022)	0.798 (\pm 0.019)	0.792 (\pm 0.022)	0.803 (\pm 0.023)	0.797 (\pm 0.021)	0.794 (\pm 0.026)	0.798 (\pm 0.025)
	SVC	0.599 (\pm 0.025)	0.786 (\pm 0.006)	0.753 (\pm 0.015)	0.781 (\pm 0.009)	0.776 (\pm 0.012)	0.795 (\pm 0.009)	0.782 (\pm 0.012)	0.793 (\pm 0.006)	0.785 (\pm 0.012)
	RandomForestClassifier	0.576 (\pm 0.036)	0.762 (\pm 0.026)	0.733 (\pm 0.026)	0.737 (\pm 0.030)	0.737 (\pm 0.035)	0.738 (\pm 0.036)	0.737 (\pm 0.042)	0.733 (\pm 0.030)	0.736 (\pm 0.033)
	XGBClassifier	0.624 (\pm 0.018)	0.791 (\pm 0.007)	0.780 (\pm 0.010)	0.800 (\pm 0.004)	0.793 (\pm 0.007)	0.798 (\pm 0.007)	0.800 (\pm 0.005)	0.797 (\pm 0.010)	0.799 (\pm 0.012)
Average precision	LogisticRegression	0.427 (\pm 0.032)	0.626 (\pm 0.038)	0.597 (\pm 0.049)	0.634 (\pm 0.031)	0.621 (\pm 0.051)	0.638 (\pm 0.023)	0.635 (\pm 0.035)	0.631 (\pm 0.031)	0.633 (\pm 0.030)
	LinearSVC	0.425 (\pm 0.032)	0.624 (\pm 0.037)	0.596 (\pm 0.028)	0.624 (\pm 0.027)	0.622 (\pm 0.032)	0.628 (\pm 0.031)	0.627 (\pm 0.031)	0.624 (\pm 0.037)	0.626 (\pm 0.038)
	SVC	0.376 (\pm 0.029)	0.623 (\pm 0.013)	0.566 (\pm 0.024)	0.608 (\pm 0.019)	0.601 (\pm 0.024)	0.631 (\pm 0.017)	0.610 (\pm 0.023)	0.628 (\pm 0.012)	0.614 (\pm 0.022)
	RandomForestClassifier	0.343 (\pm 0.030)	0.584 (\pm 0.036)	0.502 (\pm 0.028)	0.526 (\pm 0.045)	0.534 (\pm 0.046)	0.524 (\pm 0.050)	0.545 (\pm 0.057)	0.518 (\pm 0.044)	0.536 (\pm 0.050)
	XGBClassifier	0.400 (\pm 0.021)	0.645 (\pm 0.021)	0.567 (\pm 0.028)	0.625 (\pm 0.010)	0.620 (\pm 0.010)	0.640 (\pm 0.014)	0.641 (\pm 0.007)	0.635 (\pm 0.040)	0.646 (\pm 0.017)
Accuracy	LogisticRegression	0.553 (\pm 0.032)	0.612 (\pm 0.058)	0.673 (\pm 0.052)	0.683 (\pm 0.081)	0.637 (\pm 0.100)	0.639 (\pm 0.063)	0.653 (\pm 0.094)	0.631 (\pm 0.058)	0.637 (\pm 0.087)
	LinearSVC	0.569 (\pm 0.036)	0.665 (\pm 0.050)	0.687 (\pm 0.030)	0.719 (\pm 0.037)	0.699 (\pm 0.040)	0.725 (\pm 0.049)	0.706 (\pm 0.043)	0.686 (\pm 0.054)	0.708 (\pm 0.052)
	SVC	0.556 (\pm 0.071)	0.832 (\pm 0.018)	0.699 (\pm 0.023)	0.732 (\pm 0.031)	0.714 (\pm 0.038)	0.787 (\pm 0.029)	0.737 (\pm 0.038)	0.838 (\pm 0.000)	0.761 (\pm 0.047)
	RandomForestClassifier	0.723 (\pm 0.007)	0.766 (\pm 0.029)	0.742 (\pm 0.011)	0.745 (\pm 0.015)	0.751 (\pm 0.019)	0.745 (\pm 0.017)	0.753 (\pm 0.017)	0.752 (\pm 0.017)	0.755 (\pm 0.018)
	XGBClassifier	0.564 (\pm 0.065)	0.698 (\pm 0.064)	0.743 (\pm 0.014)	0.769 (\pm 0.012)	0.758 (\pm 0.021)	0.763 (\pm 0.010)	0.759 (\pm 0.016)	0.756 (\pm 0.011)	0.757 (\pm 0.017)
Balanced accuracy	LogisticRegression	0.594 (\pm 0.021)	0.676 (\pm 0.026)	0.705 (\pm 0.034)	0.719 (\pm 0.041)	0.692 (\pm 0.049)	0.698 (\pm 0.030)	0.702 (\pm 0.044)	0.687 (\pm 0.028)	0.694 (\pm 0.039)
	LinearSVC	0.608 (\pm 0.029)	0.698 (\pm 0.028)	0.716 (\pm 0.020)	0.736 (\pm 0.023)	0.721 (\pm 0.023)	0.738 (\pm 0.032)	0.725 (\pm 0.025)	0.711 (\pm 0.031)	0.727 (\pm 0.032)
	SVC	0.556 (\pm 0.031)	0.780 (\pm 0.010)	0.700 (\pm 0.012)	0.729 (\pm 0.018)	0.720 (\pm 0.022)	0.761 (\pm 0.016)	0.734 (\pm 0.023)	0.783 (\pm 0.000)	0.745 (\pm 0.026)
	RandomForestClassifier	0.523 (\pm 0.028)	0.670 (\pm 0.065)	0.585 (\pm 0.072)	0.564 (\pm 0.065)	0.590 (\pm 0.075)	0.560 (\pm 0.058)	0.591 (\pm 0.073)	0.587 (\pm 0.049)	0.591 (\pm 0.066)
	XGBClassifier	0.581 (\pm 0.022)	0.715 (\pm 0.031)	0.749 (\pm 0.008)	0.763 (\pm 0.006)	0.751 (\pm 0.013)	0.758 (\pm 0.006)	0.756 (\pm 0.008)	0.756 (\pm 0.007)	0.752 (\pm 0.010)
Sensitivity	LogisticRegression	0.685 (\pm 0.053)	0.818 (\pm 0.052)	0.777 (\pm 0.029)	0.799 (\pm 0.050)	0.814 (\pm 0.072)	0.830 (\pm 0.047)	0.812 (\pm 0.071)	0.810 (\pm 0.045)	0.820 (\pm 0.069)
	LinearSVC	0.697 (\pm 0.056)	0.772 (\pm 0.029)	0.780 (\pm 0.030)	0.774 (\pm 0.025)	0.768 (\pm 0.029)	0.766 (\pm 0.022)	0.768 (\pm 0.027)	0.767 (\pm 0.030)	0.768 (\pm 0.026)
	SVC	0.555 (\pm 0.190)	0.665 (\pm 0.007)	0.702 (\pm 0.024)	0.723 (\pm 0.016)	0.733 (\pm 0.020)	0.703 (\pm 0.018)	0.727 (\pm 0.016)	0.663 (\pm 0.001)	0.710 (\pm 0.022)
	RandomForestClassifier	0.079 (\pm 0.081)	0.459 (\pm 0.176)	0.237 (\pm 0.201)	0.165 (\pm 0.171)	0.234 (\pm 0.191)	0.152 (\pm 0.142)	0.233 (\pm 0.191)	0.223 (\pm 0.115)	0.230 (\pm 0.168)
	XGBClassifier	0.619 (\pm 0.085)	0.753 (\pm 0.041)	0.763 (\pm 0.009)	0.751 (\pm 0.007)	0.736 (\pm 0.013)	0.747 (\pm 0.006)	0.748 (\pm 0.015)	0.757 (\pm 0.006)	0.741 (\pm 0.017)
Specificity	LogisticRegression	0.503 (\pm 0.054)	0.534 (\pm 0.094)	0.633 (\pm 0.073)	0.640 (\pm 0.123)	0.570 (\pm 0.156)	0.567 (\pm 0.099)	0.593 (\pm 0.149)	0.563 (\pm 0.091)	0.568 (\pm 0.138)
	LinearSVC	0.520 (\pm 0.054)	0.625 (\pm 0.074)	0.651 (\pm 0.043)	0.698 (\pm 0.053)	0.673 (\pm 0.059)	0.710 (\pm 0.069)	0.682 (\pm 0.063)	0.655 (\pm 0.078)	0.685 (\pm 0.074)
	SVC	0.557 (\pm 0.159)	0.896 (\pm 0.027)	0.697 (\pm 0.036)	0.736 (\pm 0.044)	0.707 (\pm 0.055)	0.819 (\pm 0.044)	0.741 (\pm 0.054)	0.904 (\pm 0.001)	0.780 (\pm 0.068)
	RandomForestClassifier	0.966 (\pm 0.030)	0.882 (\pm 0.070)	0.932 (\pm 0.065)	0.964 (\pm 0.049)	0.945 (\pm 0.051)	0.969 (\pm 0.032)	0.949 (\pm 0.053)	0.952 (\pm 0.025)	0.953 (\pm 0.044)
	XGBClassifier	0.543 (\pm 0.115)	0.678 (\pm 0.099)	0.736 (\pm 0.021)	0.775 (\pm 0.018)	0.767 (\pm 0.029)	0.769 (\pm 0.014)	0.764 (\pm 0.025)	0.756 (\pm 0.015)	0.763 (\pm 0.026)

Table A.6: Predictive performance on DIGPD of the five machine learning algorithms with the nine reduction approaches with 10 repetitions of the nest cross-validation. Mean (standard deviation) over the 10 models are reported.

Bibliography

J. Eric Ahlskog. Does vigorous exercise have a neuroprotective effect in Parkinson disease? *Neurology*, 77(3):288–294, July 2011. ISSN 1526-632X. doi: 10.1212/WNL.0b013e318225ab66.

American Psychiatric Association. *Diagnostic and Statistical Manual of Mental Disorders, 5th Edition*. American Psychiatric Pub, May 2013. ISBN 978-0-89042-557-2.

Elise Anderson and John Nutt. The long-duration response to levodopa: Phenomenology, potential mechanisms and clinical implications. *Parkinsonism & Related Disorders*, 17(8):587–592, September 2011. ISSN 1353-8020. doi: 10.1016/j.parkreldis.2011.03.014. URL <http://www.sciencedirect.com/science/article/pii/S1353802011000873>.

Camila Catherine Aquino and Susan H. Fox. Clinical spectrum of levodopa-induced complications. *Movement Disorders*, 30(1):80–89, 2015. ISSN 1531-8257. doi: 10.1002/mds.26125. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/mds.26125>. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/mds.26125>.

C. Ardouin, I. Chéreau, P. M. Llorca, E. Lhommée, F. Durif, P. Pollak, and P. Krack. Évaluation des troubles comportementaux hyper- et hypodopaminergiques dans la maladie de Parkinson. *Revue Neurologique*, 165(11):845–856, November 2009. ISSN 0035-3787. doi: 10.1016/j.neurol.2009.06.003. URL <http://www.sciencedirect.com/science/article/pii/S0035378709003439>.

M. Auyeung, T. H. Tsoi, W. K. Tang, C. M. Cheung, C. N. Lee, R. Li, and Eric Yeung. Impulse control disorders in Chinese Parkinson’s disease patients: the effect of ergot derived dopamine agonist. *Parkinsonism & Related Disorders*, 17(8):635–637, September 2011. ISSN 1873-5126. doi: 10.1016/j.parkreldis.2011.06.001.

Fahd Baig, Mark J. Kelly, Michael A. Lawton, Claudio Ruffmann, Michal Rolinski, Johannes C. Klein, Thomas Barber, Christine Lo, Yoav Ben-Shlomo, David Okai, and Michele T. Hu. Impulse control disorders in Parkinson disease and RBD: A longitudinal study of severity. *Neurology*, 93(7):e675–e687, 2019. ISSN 1526-632X. doi: 10.1212/WNL.0000000000007942.

- Sara Balduzzi, Gerta Rücker, and Guido Schwarzer. How to perform a meta-analysis with R: a practical tutorial. *Evidence-Based Mental Health*, 22(4):153–160, November 2019. ISSN 1468-960X. doi: 10.1136/ebmental-2019-300117.
- Roberta Balestrino and Pablo Martinez-Martin. Neuropsychiatric symptoms, behavioural disorders, and quality of life in Parkinson’s disease. *Journal of the Neurological Sciences*, 373:173–178, February 2017. ISSN 1878-5883. doi: 10.1016/j.jns.2016.12.060.
- Jesse Bastiaens, Benjamin J. Dorfman, Paul J. Christos, and Melissa J. Nirenberg. Prospective cohort study of impulse control disorders in Parkinson’s disease. *Movement Disorders: Official Journal of the Movement Disorder Society*, 28(3):327–333, March 2013. ISSN 1531-8257. doi: 10.1002/mds.25291.
- V. Biouesse, B. C. Skibell, R. L. Watts, D. N. Loupe, C. Drews-Botsch, and N. J. Newman. Ophthalmologic features of Parkinson’s disease. *Neurology*, 62(2):177–180, January 2004. ISSN 1526-632X. doi: 10.1212/01.wnl.0000103444.45882.d8.
- D. W. Black and J. E. Grant. *DSM-5 Guidebook: The Essential Companion to the Diagnostic and Statistical Manual of Mental Disorders, Fifth Edition*. American Psychiatric Pub, February 2014. ISBN 978-1-58562-465-2. Google-Books-ID: IKeTAwAAQBAJ.
- Donald W. Black. A review of compulsive buying disorder. *World Psychiatry*, 6(1):14–18, February 2007. ISSN 1723-8617. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1805733/>.
- Bernhard E. Boser, Isabelle M. Guyon, and Vladimir N. Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the fifth annual workshop on Computational learning theory*, COLT ’92, pages 144–152, Pittsburgh, Pennsylvania, USA, July 1992. Association for Computing Machinery. ISBN 978-0-89791-497-0. doi: 10.1145/130385.130401. URL <https://doi.org/10.1145/130385.130401>.
- Leo Breiman. Bias, Variance , AND Arcing Classifiers. Technical Report 460, Statistics Department University of California, Berkeley, 1996.
- Leo Breiman. Arcing the Edge. Technical Report 486, Statistics Department University of California, Berkeley, 1997.
- Leo Breiman. Random Forests. *Machine Learning*, 45(1):5–32, October 2001. ISSN 1573-0565. doi: 10.1023/A:1010933404324. URL <https://doi.org/10.1023/A:1010933404324>.
- Leo Breiman, Jerome Friedman, Charles J. Stone, and R. A. Olshen. *Classification and Regression Trees*. Taylor & Francis, January 1984. ISBN 978-0-412-04841-8.

- Martijn P. G. Broen, Nadia E. Narayen, Mark L. Kuijf, Nadeeka N. W. Dissanayaka, and Albert F. G. Leentjens. Prevalence of anxiety in Parkinson's disease: A systematic review and meta-analysis. *Movement Disorders: Official Journal of the Movement Disorder Society*, 31(8):1125–1133, 2016. ISSN 1531-8257. doi: 10.1002/mds.26643.
- Annalisa Buniello, Jacqueline A. L. MacArthur, Maria Cerezo, Laura W. Harris, James Hayhurst, Cinzia Malangone, Aoife McMahon, Joannella Morales, Edward Mountjoy, Elliot Sollis, Daniel Suveges, Olga Vrousseau, Patricia L. Whetzel, Ridwan Amode, Jose A. Guillen, Harpreet S. Riat, Stephen J. Trevanion, Peggy Hall, Heather Jenkins, Paul Flicek, Tony Burdett, Lucia A. Hindorff, Fiona Cunningham, and Helen Parkinson. The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Research*, 47(D1):D1005–D1012, 2019. ISSN 1362-4962. doi: 10.1093/nar/gky1120.
- M. B. Callesen, D. Weintraub, M. F. Damholdt, and A. Møller. Impulsive and compulsive behaviors among Danish patients with Parkinson's disease: prevalence, depression, and personality. *Parkinsonism & Related Disorders*, 20(1):22–26, January 2014. ISSN 1873-5126. doi: 10.1016/j.parkreldis.2013.09.006.
- Mette Buhl Callesen and Malene Flensburg Damholdt. Phenomenology and gender characteristics of hobbyism and punning in Parkinson's disease: A self-report study. *Basal Ganglia*, 9:1–6, August 2017. ISSN 2210-5336. doi: 10.1016/j.baga.2017.06.002. URL <http://www.sciencedirect.com/science/article/pii/S2210533616300600>.
- D. B. Calne, B. J. Snow, and C. Lee. Criteria for diagnosing Parkinson's disease. *Annals of Neurology*, 32 Suppl:S125–127, 1992. ISSN 0364-5134. doi: 10.1002/ana.410320721.
- Arvid Carlsson, Margit Lindqvist, and Tor Magnusson. 3,4-Dihydroxyphenylalanine and 5-Hydroxytryptophan as Reserpine Antagonists. *Nature*, 180(4596):1200–1200, November 1957. ISSN 1476-4687. doi: 10.1038/1801200a0. URL <https://www.nature.com/articles/1801200a0>. Number: 4596 Publisher: Nature Publishing Group.
- Jodi Cartoon and Jothi Ramalingam. Dopamine dysregulation syndrome in non-Parkinson's disease patients: a systematic review. *Australasian Psychiatry: Bulletin of Royal Australian and New Zealand College of Psychiatrists*, 27(5):456–461, October 2019. ISSN 1440-1665. doi: 10.1177/1039856219839476.
- Xochitl Helga Castro-Martínez, Pedro J. García-Ruiz, Carlos Martínez-García, Juan Carlos Martínez-Castrillo, Lydia Vela, Marina Mata, Irene Martínez-Torres, Cici Feliz-Feliz, Francesc Palau, and Janet Hoenicka. Behavioral addictions in early-onset Parkinson disease are associated with DRD3 variants. *Parkinsonism & Related Disorders*, 49:100–103, 2018. ISSN 1873-5126. doi: 10.1016/j.parkreldis.2018.01.010.

- Lama M. Chahine, Amy W. Amara, and Aleksandar Videnovic. A systematic review of the literature on disorders of sleep and wakefulness in Parkinson's disease from 2005 to 2015. *Sleep Medicine Reviews*, 35:33–50, 2017. ISSN 1532-2955. doi: 10.1016/j.smr.2016.08.001.
- Samuel R. Chamberlain and Jon E. Grant. Minnesota Impulse Disorders Interview (MIDI): Validation of a structured diagnostic clinical interview for impulse control disorders in an enriched community sample. *Psychiatry Research*, 265:279–283, July 2018. ISSN 0165-1781. doi: 10.1016/j.psychres.2018.05.006. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5985960/>.
- Christopher C. Chang, Carson C. Chow, Laurent CAM Tellier, Shashaank Vattikuti, Shaun M. Purcell, and James J. Lee. Second-generation PLINK: rising to the challenge of larger and richer datasets. *GigaScience*, 4(1), December 2015. doi: 10.1186/s13742-015-0047-8. URL <https://academic.oup.com/gigascience/article/4/1/s13742-015-0047-8/2707533>. Publisher: Oxford Academic.
- Fanny Charbonnier-Beaupel, Marion Malerbi, Cristina Alcacer, Khadija Tahiri, Wasila Carpentier, Chuansong Wang, Matthew During, Desheng Xu, Paul F. Worley, Jean-Antoine Girault, Denis Hervé, and Jean-Christophe Corvol. Gene Expression Analyses Identify Narp Contribution in the Development of l-DOPA-Induced Dyskinesia. *Journal of Neuroscience*, 35(1):96–111, January 2015. ISSN 0270-6474, 1529-2401. doi: 10.1523/JNEUROSCI.5231-13.2015. URL <https://www.jneurosci.org/content/35/1/96>. Publisher: Society for Neuroscience Section: Articles.
- K. Ray Chaudhuri, Daniel G. Healy, Anthony H. V. Schapira, and National Institute for Clinical Excellence. Non-motor symptoms of Parkinson's disease: diagnosis and management. *The Lancet. Neurology*, 5(3):235–245, March 2006. ISSN 1474-4422. doi: 10.1016/S1474-4422(06)70373-8.
- Tianqi Chen and Carlos Guestrin. XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 785–794, San Francisco California USA, August 2016. ACM. ISBN 978-1-4503-4232-2. doi: 10.1145/2939672.2939785. URL <https://dl.acm.org/doi/10.1145/2939672.2939785>.
- Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning Phrase Representations using RNN EncoderDecoder for Statistical Machine Translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.3115/v1/D14-1179. URL <https://www.aclweb.org/anthology/D14-1179>.

- Ji Hyun Choi, Jee Young Lee, Jin Whan Cho, Seong Beom Ko, Tae Beom Ahn, Sang Jin Kim, Sang Myung Cheon, Joong Seok Kim, Yoon Joong Kim, Hyeo Il Ma, Jong Sam Baik, Phil Hyu Lee, Sun Ju Chung, Jong Min Kim, In Uk Song, Han Joon Kim, Young Hee Sung, Do Young Kwon, Jae Hyeok Lee, Ji Young Kim, Ji Sun Kim, Ji Young Yun, Hee Jin Kim, Jin Yong Hong, Mi Jung Kim, Jinyoung Youn, Ji Seon Kim, Eung Seok Oh, Hui Jun Yang, Won Tae Yoon, Sooyeoun You, Kyum Yil Kwon, Hyung Eun Park, Su Yun Lee, Younsoo Kim, Hee Tae Kim, and Mee Young Park. Validation of the Korean Version of the Questionnaire for Impulsive-Compulsive Disorders in Parkinson's Disease Rating Scale. *Journal of Clinical Neurology (Seoul, Korea)*, 16(2):245–253, April 2020. ISSN 1738-6586. doi: 10.3988/jcn.2020.16.2.245.
- G. A. Christenson, R. J. Faber, M. de Zwaan, N. C. Raymond, S. M. Specker, M. D. Ekern, T. B. Mackenzie, R. D. Crosby, S. J. Crow, and E. D. Eckert. Compulsive buying: descriptive characteristics and psychiatric comorbidity. *The Journal of Clinical Psychiatry*, 55(1):5–11, January 1994. ISSN 0160-6689.
- Junyoung Chung, Caglar Gulcehre, Kyunghyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *NIPS 2014 Workshop on Deep Learning, December 2014*, 2014. URL <https://nyuscholars.nyu.edu/en/publications/empirical-evaluation-of-gated-recurrent-neural-networks-on-sequen>.
- Roberto Cilia, Chiara Siri, Margherita Canesi, Anna Lena Zecchinelli, Danilo De Gaspari, Francesca Natuzzi, Silvana Tesei, Nicoletta Meucci, Claudio Bruno Mariani, Giorgio Sacilotto, Michela Zini, Claudio Ruffmann, and Gianni Pezzoli. Dopamine dysregulation syndrome in Parkinson's disease: from clinical and neuropsychological characterisation to management and long-term outcome. *Journal of Neurology, Neurosurgery, and Psychiatry*, 85(3):311–318, March 2014. ISSN 1468-330X. doi: 10.1136/jnnp-2012-303988.
- Barbara S. Connolly and Anthony E. Lang. Pharmacological treatment of Parkinson disease: a review. *JAMA*, 311(16):1670–1683, April 2014. ISSN 1538-3598. doi: 10.1001/jama.2014.3654.
- Gregory Cooper, Gerald Eichhorn, and Robert L. Rodnitzky. Parkinson's Disease. In *Neuroscience in Medicine*, pages 508–512. Humana Press, 3 edition, 2008. ISBN 978-1-60327-455-5.
- Florence Cormier-Dequaire, Samir Bekadar, Mathieu Anheim, Said Lebbah, Antoine Pelissolo, Paul Krack, Lucette Lacomblez, Eugénie Lhommée, Anna Castrioto, Jean-Philippe Azulay, Luc Defebvre, Alexandre Kreisler, Franck Durif, Ana Marques-Raquel, Christine Brefel-Courbon, David Grabli, Emmanuel Roze, Pierre-Michel Llorca, Fabienne Ory-Magne, Isabelle Benatru, Solene Ansquer, David Maltête,

- Melissa Tir, Pierre Krystkowiak, Christine Tranchant, Ouhaïd Lagha-Boukbiza, Bénédicte Lebrun-Vignes, Graziella Mangone, Marie Vidailhet, Fanny Charbonnier-Beaupel, Olivier Rascol, Suzanne Lesage, Alexis Brice, Sophie Tezenas du Montcel, Jean-Christophe Corvol, and BADGE-PD study group. Suggestive association between OPRM1 and impulse control disorders in Parkinson's disease. *Movement Disorders: Official Journal of the Movement Disorder Society*, 33(12):1878–1886, 2018. ISSN 1531-8257. doi: 10.1002/mds.27519.
- Jason R. Cornelius, Maja Tippmann-Peikert, Nancy L. Slocumb, Courtney F. Frerichs, and Michael H. Silber. Impulse control disorders with the use of dopaminergic agents in restless legs syndrome: a case-control study. *Sleep*, 33(1):81–87, January 2010. ISSN 0161-8105.
- Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, September 1995. ISSN 1573-0565. doi: 10.1007/BF00994018. URL <https://doi.org/10.1007/BF00994018>.
- Jean-Christophe Corvol, Jean-Baptiste Anzouan-Kacou, Elodie Fauveau, Anne-Marie Bonnet, Bénédicte Lebrun-Vignes, Camille Girault, Yves Agid, Philippe Lechat, Richard Isnard, and Lucette Lacomblez. Heart valve regurgitation, pergolide use, and parkinson disease: an observational study and meta-analysis. *Archives of Neurology*, 64(12):1721–1726, December 2007. ISSN 0003-9942. doi: 10.1001/archneur.64.12.1721.
- Jean-Christophe Corvol, Fanny Artaud, Florence Cormier-Dequaire, Olivier Rascol, Franck Durif, Pascal Derkinderen, Ana-Raquel Marques, Frédéric Bourdain, Jean-Philippe Brandel, Fernando Pico, Lucette Lacomblez, Cecilia Bonnet, Christine Brefel-Courbon, Fabienne Ory-Magne, David Grabli, Stephan Klebe, Graziella Mangone, Hana You, Valérie Mesnage, Pei-Chen Lee, Alexis Brice, Marie Vidailhet, Alexis Elbaz, and DIGPD Study Group. Longitudinal analysis of impulse control disorders in Parkinson disease. *Neurology*, 91(3):e189–e201, July 2018. ISSN 1526-632X. doi: 10.1212/WNL.0000000000005816.
- G. C. Cotzias, M. H. Van Woert, and L. M. Schiffer. Aromatic amino acids and modification of parkinsonism. *The New England Journal of Medicine*, 276(7):374–379, February 1967. ISSN 0028-4793. doi: 10.1056/NEJM196702162760703.
- Sayantan Das, Lukas Forer, Sebastian Schönherr, Carlo Sidore, Adam E. Locke, Alan Kwong, Scott I. Vrieze, Emily Y. Chew, Shawn Levy, Matt McGue, David Schlessinger, Dwight Stambolian, Po-Ru Loh, William G. Iacono, Anand Swaroop, Laura J. Scott, Francesco Cucca, Florian Kronenberg, Michael Boehnke, Gonçalo R. Abecasis, and Christian Fuchsberger. Next-generation genotype imputation service

- and methods. *Nature Genetics*, 48(10):1284–1287, 2016. ISSN 1546-1718. doi: 10.1038/ng.3656.
- C. A. Davie. A review of Parkinson’s disease. *British Medical Bulletin*, 86:109–127, 2008. ISSN 1471-8391. doi: 10.1093/bmb/ldn013.
- Jesse Davis and Mark Goadrich. The relationship between Precision-Recall and ROC curves. In *Proceedings of the 23rd international conference on Machine learning, ICML ’06*, pages 233–240, Pittsburgh, Pennsylvania, USA, June 2006. Association for Computing Machinery. ISBN 978-1-59593-383-6. doi: 10.1145/1143844.1143874. URL <https://doi.org/10.1145/1143844.1143874>.
- Patricia de la Riva, Kara Smith, Sharon X. Xie, and Daniel Weintraub. Course of psychiatric symptoms and global cognition in early Parkinson disease. *Neurology*, 83(12): 1096–1103, September 2014. ISSN 1526-632X. doi: 10.1212/WNL.0000000000000801.
- Bernardo Dell’Osso, A. Carlo Altamura, Andrea Allen, Donatella Marazziti, and Eric Hollander. Epidemiologic and clinical updates on impulse control disorders: a critical review. *European Archives of Psychiatry and Clinical Neuroscience*, 256(8):464–475, December 2006. ISSN 0940-1334. doi: 10.1007/s00406-006-0668-0.
- E. R. DeLong, D. M. DeLong, and D. L. Clarke-Pearson. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. *Biometrics*, 44(3):837–845, September 1988. ISSN 0006-341X.
- Ditte Demontis, Raymond K. Walters, Joanna Martin, Manuel Mattheisen, Thomas D. Als, Esben Agerbo, Gísli Baldursson, Rich Belliveau, Jonas Bybjerg-Grauholm, Marie Bækvad-Hansen, Felecia Cerrato, Kimberly Chambert, Claire Churchhouse, Ashley Dumont, Nicholas Eriksson, Michael Gandal, Jacqueline I. Goldstein, Katrina L. Grasby, Jakob Grove, Olafur O. Gudmundsson, Christine S. Hansen, Mads Engel Hauberg, Mads V. Hollegaard, Daniel P. Howrigan, Hailiang Huang, Julian B. Maller, Alicia R. Martin, Nicholas G. Martin, Jennifer Moran, Jonatan Pallesen, Duncan S. Palmer, Carsten Bøcker Pedersen, Marianne Giørtz Pedersen, Timothy Poterba, Jesper Buchhave Poulsen, Stephan Ripke, Elise B. Robinson, F. Kyle Satterstrom, Hreinn Stefansson, Christine Stevens, Patrick Turley, G. Bragi Walters, Hyejung Won, Margaret J. Wright, ADHD Working Group of the Psychiatric Genomics Consortium (PGC), Early Lifecourse & Genetic Epidemiology (EAGLE) Consortium, 23andMe Research Team, Ole A. Andreassen, Philip Asherson, Christie L. Burton, Dorret I. Boomsma, Bru Cormand, Søren Dalsgaard, Barbara Franke, Joel Gelernter, Daniel Geschwind, Hakon Hakonarson, Jan Haavik, Henry R. Kranzler, Jonna Kuntsi, Kate Langley, Klaus-Peter Lesch, Christel Middeldorp, Andreas Reif, Luis Augusto Rohde, Panos Roussos, Russell Schachar, Pamela Sklar, Edmund J. S.

- Sonuga-Barke, Patrick F. Sullivan, Anita Thapar, Joyce Y. Tung, Irwin D. Waldman, Sarah E. Medland, Kari Stefansson, Merete Nordentoft, David M. Hougaard, Thomas Werge, Ole Mors, Preben Bo Mortensen, Mark J. Daly, Stephen V. Faraone, Anders D. Børglum, and Benjamin M. Neale. Discovery of the first genome-wide significant risk loci for attention deficit/hyperactivity disorder. *Nature Genetics*, 51(1):63–75, 2019. ISSN 1546-1718. doi: 10.1038/s41588-018-0269-7.
- Dennis W. Dickson. Neuropathology of Parkinson disease. *Parkinsonism & Related Disorders*, 46 Suppl 1:S30–S33, January 2018. ISSN 1873-5126. doi: 10.1016/j.parkreldis.2017.07.033.
- E. Driver-Dunckley, J. Samanta, and M. Stacy. Pathological gambling associated with dopamine agonist therapy in Parkinson’s disease. *Neurology*, 61(3):422–423, August 2003. ISSN 1526-632X. doi: 10.1212/01.wnl.0000076478.45005.ec.
- Frank Dudbridge. Power and Predictive Accuracy of Polygenic Risk Scores. *PLOS Genetics*, 9(3):e1003348, March 2013. ISSN 1553-7404. doi: 10.1371/journal.pgen.1003348. URL <https://journals.plos.org/plosgenetics/article?id=10.1371/journal.pgen.1003348>. Publisher: Public Library of Science.
- H. Ehringer and O. Hornykiewicz. Verteilung Von Noradrenalin Und Dopamin (3-Hydroxytyramin) Im Gehirn Des Menschen Und Ihr Verhalten Bei Erkrankungen Des Extrapiramidalen Systems. *Klinische Wochenschrift*, 38(24):1236–1239, December 1960. ISSN 1432-1440. doi: 10.1007/BF01485901. URL <https://doi.org/10.1007/BF01485901>.
- H. El Otmani, L. Raji, B. El Moutaouakil, M. A. Rafai, and I. Slassi. Punding sévère au cours d’une maladie de Parkinson. *L’Encéphale*, 41(2):190–193, April 2015. ISSN 0013-7006. doi: 10.1016/j.encep.2013.03.013. URL <http://www.sciencedirect.com/science/article/pii/S0013700613001656>.
- Aleksander H. Erga, Ingvid Dalen, Anastasia Ushakova, Janete Chung, Charalampos Tzoulis, Ole Bjørn Tysnes, Guido Alves, Kenn Freddy Pedersen, and Jodi Maple-Grødem. Dopaminergic and Opioid Pathways Associated with Impulse Control Disorders in Parkinsons Disease. *Frontiers in Neurology*, 9, 2018. ISSN 1664-2295. doi: 10.3389/fneur.2018.00109. URL <https://www.frontiersin.org/articles/10.3389/fneur.2018.00109/full>. Publisher: Frontiers.
- Andrew H. Evans, Regina Katzenschlager, Dominic Paviour, John D. O’Sullivan, Silke Appel, Andrew D. Lawrence, and Andrew J. Lees. Punding in Parkinson’s disease: Its relation to the dopamine dysregulation syndrome. *Movement Disorders*, 19(4):397–405, 2004. ISSN 1531-8257. doi: 10.1002/mds.20045. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/mds.20045>. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/mds.20045>.

- Andrew H. Evans, David Okai, Daniel Weintraub, Shen-Yang Lim, Sean S. O’Sullivan, Valerie Voon, Paul Krack, Cristina Sampaio, Bart Post, Albert F. G. Leentjens, Pablo Martinez-Martin, Glenn T. Stebbins, Christopher G. Goetz, Anette Schrag, and Members of the International Parkinson and Movement Disorder Society (IPMDS) Rating Scales Review Committee. Scales to assess impulsive and compulsive behaviors in Parkinson’s disease: Critique and recommendations. *Movement Disorders: Official Journal of the Movement Disorder Society*, 34(6):791–798, 2019. ISSN 1531-8257. doi: 10.1002/mds.27689.
- Stanley Fahn. The history of dopamine and levodopa in the treatment of Parkinson’s disease. *Movement Disorders*, 23(S3):S497–S508, 2008. ISSN 1531-8257. doi: 10.1002/mds.22028. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/mds.22028>. __eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/mds.22028>.
- M. L. Fantini, L. Macedo, M. Zibetti, M. Sarchioto, T. Vidal, B. Pereira, A. Marques, B. Debilly, P. Derost, M. Ulla, N. Vitello, A. Cicolin, L. Lopiano, and F. Durif. Increased risk of impulse control symptoms in Parkinson’s disease with REM sleep behaviour disorder. *Journal of Neurology, Neurosurgery, and Psychiatry*, 86(2):174–179, February 2015. ISSN 1468-330X. doi: 10.1136/jnnp-2014-307904.
- Maria Livia Fantini, Michela Figorilli, Isabelle Arnulf, Maurizio Zibetti, Bruno Pereira, Patricia Beudin, Monica Puligheddu, Florence Cormier-Dequaire, Lucette Lacomblez, Eve Benchetrit, Jean Christophe Corvol, Alessandro Cicolin, Leonardo Lopiano, Ana Marques, and Franck Durif. Sleep and REM sleep behaviour disorder in Parkinson’s disease with impulse control disorder. *Journal of Neurology, Neurosurgery, and Psychiatry*, 89(3):305–310, 2018. ISSN 1468-330X. doi: 10.1136/jnnp-2017-316576.
- Maria Livia Fantini, Franck Durif, and Ana Marques. Impulse Control Disorders in REM Sleep Behavior Disorder. *Current Treatment Options in Neurology*, 21(5):23, April 2019. ISSN 1092-8480. doi: 10.1007/s11940-019-0564-3.
- Maria Livia Fantini, Janel Fedler, Bruno Pereira, Daniel Weintraub, Ana-Raquel Marques, and Franck Durif. Is Rapid Eye Movement Sleep Behavior Disorder a Risk Factor for Impulse Control Disorder in Parkinson Disease? *Annals of Neurology*, 88(4):759–770, October 2020. ISSN 1531-8249. doi: 10.1002/ana.25798.
- Alfonso Fasano, Naomi P. Visanji, Louis W. C. Liu, Antony E. Lang, and Ronald F. Pfeiffer. Gastrointestinal dysfunction in Parkinson’s disease. *The Lancet. Neurology*, 14(6):625–639, June 2015. ISSN 1474-4465. doi: 10.1016/S1474-4422(15)00007-1.
- Leslie J. Findley. The economic impact of Parkinson’s disease. *Parkinsonism & Related Disorders*, 13 Suppl:S8–S12, September 2007. ISSN 1353-8020. doi: 10.1016/j.parkreldis.2007.06.003.

- Naomi A. Fineberg, Marc N. Potenza, Samuel R. Chamberlain, Heather A. Berlin, Lara Menzies, Antoine Bechara, Barbara J. Sahakian, Trevor W. Robbins, Edward T. Bullmore, and Eric Hollander. Probing compulsive and impulsive behaviors, from animal models to endophenotypes: a narrative review. *Neuropsychopharmacology: Official Publication of the American College of Neuropsychopharmacology*, 35(3):591–604, February 2010. ISSN 1740-634X. doi: 10.1038/npp.2009.185.
- Peter Flach and Meelis Kull. Precision-Recall-Gain Curves: PR Analysis Done Right. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 838–846. Curran Associates, Inc., 2015. URL <http://papers.nips.cc/paper/5867-precision-recall-gain-curves-pr-analysis-done-right.pdf>.
- Leonardo F. Fontenelle, Mauro V. Mendlowicz, and Marcio Versiani. Impulse control disorders in patients with obsessive-compulsive disorder. *Psychiatry and Clinical Neurosciences*, 59(1):30–37, February 2005. ISSN 1323-1316. doi: 10.1111/j.1440-1819.2005.01328.x.
- Leonardo F. Fontenelle, Sanne Oostermeijer, Ben J. Harrison, Christos Pantelis, and Murat Yücel. Obsessive-compulsive disorder, impulse control disorders and drug addiction: common features and potential treatments. *Drugs*, 71(7):827–840, May 2011. ISSN 1179-1950. doi: 10.2165/11591790-000000000-00000.
- Susan H. Fox, Regina Katzenschlager, Shen-Yang Lim, Brandon Barton, Rob M. A. de Bie, Klaus Seppi, Miguel Coelho, Cristina Sampaio, and Movement Disorder Society Evidence-Based Medicine Committee. International Parkinson and movement disorder society evidence-based medicine review: Update on treatments for the motor symptoms of Parkinson’s disease. *Movement Disorders: Official Journal of the Movement Disorder Society*, 33(8):1248–1266, 2018. ISSN 1531-8257. doi: 10.1002/mds.27372.
- Jerome H. Friedman. Greedy Function Approximation: A Gradient Boosting Machine. *The Annals of Statistics*, 29(5):1189–1232, 2001. ISSN 0090-5364. URL <https://www.jstor.org/stable/2699986>. Publisher: Institute of Mathematical Statistics.
- Pedro J. Garcia-Ruiz, Juan Carlos Martinez Castrillo, Araceli Alonso-Canovas, Antonio Herranz Barcenas, Lydia Vela, Pilar Sanchez Alonso, Marina Mata, Nuria Olmedilla Gonzalez, and Ignacio Mahillo Fernandez. Impulse control disorder in patients with Parkinson’s disease under dopamine agonist therapy: a multicentre study. *Journal of Neurology, Neurosurgery, and Psychiatry*, 85(8):840–844, August 2014. ISSN 1468-330X. doi: 10.1136/jnnp-2013-306787.
- Chenjie Ge, Irene Yu-Hua Gu, Asgeir Store Jakola, and Jie Yang. Deep Learning and Multi-Sensor Fusion for Glioma Classification Using Multistream 2D Convolutional

- Networks. In *2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, pages 5894–5897, Honolulu, HI, July 2018. IEEE. ISBN 978-1-5386-3646-6. doi: 10.1109/EMBC.2018.8513556. URL <https://ieeexplore.ieee.org/document/8513556/>.
- T. Gescheidt, V. Majerová, K. Meníková, L. Duek, K. Czekóová, P. Kotková, P. Kaovský, J. Roth, and M. Bare. ID 16 Impulse control disorders in young-onset patients with Parkinsons disease: Cross-sectional study seeking associated factors with regard of personal characteristics. *Clinical Neurophysiology*, 127(3):e70, March 2016. ISSN 1388-2457. doi: 10.1016/j.clinph.2015.11.233. URL <http://www.sciencedirect.com/science/article/pii/S1388245715013383>.
- Michel Goedert, Maria Grazia Spillantini, Kelly Del Tredici, and Heiko Braak. 100 years of Lewy pathology. *Nature Reviews. Neurology*, 9(1):13–24, 2013. ISSN 1759-4766. doi: 10.1038/nrneurol.2012.242.
- Christopher G. Goetz, Barbara C. Tilley, Stephanie R. Shaftman, Glenn T. Stebbins, Stanley Fahn, Pablo Martinez-Martin, Werner Poewe, Cristina Sampaio, Matthew B. Stern, Richard Dodel, Bruno Dubois, Robert Holloway, Joseph Jankovic, Jaime Kulisevsky, Anthony E. Lang, Andrew Lees, Sue Leurgans, Peter A. LeWitt, David Nyenhuis, C. Warren Olanow, Olivier Rascol, Anette Schrag, Jeanne A. Teresi, Jacobus J. van Hilten, Nancy LaPelle, and Movement Disorder Society UPDRS Revision Task Force. Movement Disorder Society-sponsored revision of the Unified Parkinson’s Disease Rating Scale (MDS-UPDRS): scale presentation and clinimetric testing results. *Movement Disorders: Official Journal of the Movement Disorder Society*, 23(15): 2129–2170, November 2008. ISSN 1531-8257. doi: 10.1002/mds.22340.
- Jennifer G. Goldman, Samantha K. Holden, Irene Litvan, Ian McKeith, Glenn T. Stebbins, and John-Paul Taylor. Evolution of diagnostic criteria and assessments for Parkinson’s disease mild cognitive impairment. *Movement Disorders: Official Journal of the Movement Disorder Society*, 33(4):503–510, 2018. ISSN 1531-8257. doi: 10.1002/mds.27323.
- Marleide da Mota Gomes and Elias Engelhardt. Jean-Martin Charcot, father of modern neurology: an homage 120 years after his death. *Arquivos De Neuro-Psiquiatria*, 71(10):815–817, October 2013. ISSN 1678-4227. doi: 10.1590/0004-282X20130128.
- Zahra Goodarzi, Kelly J. Mrklas, Derek J. Roberts, Nathalie Jette, Tamara Pringsheim, and Jayna Holroyd-Leduc. Detecting depression in Parkinson disease: A systematic review and meta-analysis. *Neurology*, 87(4):426–437, July 2016. ISSN 1526-632X. doi: 10.1212/WNL.0000000000002898.
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. URL <http://www.deeplearningbook.org>.

- Marie Grall-Bronnec, Caroline Victorri-Vigneau, Yann Donnio, Juliette Leboucher, Morgane Rousselet, Elsa Thiabaud, Nicolas Zreika, Pascal Derkinderen, and Gaëlle Challet-Bouju. Dopamine Agonists and Impulse Control Disorders: A Complex Association. *Drug Safety*, 41(1):19–75, January 2018. ISSN 1179-1942. doi: 10.1007/s40264-017-0590-6.
- Jon E. Grant. *Impulse Control Disorders: A Clinician's Guide to Understanding and Treating Behavioral Addictions*. WW Norton and Company, New York, 2008.
- Jon E. Grant, Maria C. Mancebo, Anthony Pinto, Jane L. Eisen, and Steven A. Rasmussen. Impulse control disorders in adults with obsessive compulsive disorder. *Journal of Psychiatric Research*, 40(6):494–501, September 2006. ISSN 0022-3956. doi: 10.1016/j.jpsychires.2005.11.005.
- Jon E. Grant, Maria C. Mancebo, Jane L. Eisen, and Steven A. Rasmussen. Impulse-control disorders in children and adolescents with obsessive-compulsive disorder. *Psychiatry Research*, 175(1-2):109–113, January 2010. ISSN 0165-1781. doi: 10.1016/j.psychres.2009.04.006.
- Ute Gschwandtner, Jacqueline Aston, Susanne Renaud, and Peter Fuhr. Pathologic Gambling in Patients with Parkinson's Disease:. *Clinical Neuropharmacology*, 24(3): 170–172, May 2001. ISSN 0362-5664. doi: 10.1097/00002826-200105000-00009. URL <http://journals.lww.com/00002826-200105000-00009>.
- Jie Hao, Sai Chandra Kosaraju, Nelson Zange Tsaku, Dae Hyun Song, and Mingon Kang. PAGE-Net: Interpretable and Integrative Deep Learning for Survival Analysis Using Histopathological Images and Genomic Data. In *Biocomputing 2020*, pages 355–366. WORLD SCIENTIFIC, November 2019. ISBN 9789811215629. doi: 10.1142/9789811215636_0032. URL https://www.worldscientific.com/doi/abs/10.1142/9789811215636_0032. ZSCC: 0000001.
- Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, September 2020a. ISSN 1476-4687. doi: 10.1038/s41586-020-2649-2. URL <https://www.nature.com/articles/s41586-020-2649-2>. Number: 7825 Publisher: Nature Publishing Group.
- James P. Harris, Justin C. Burrell, Laura A. Struzyna, H. Isaac Chen, Mijail D. Serruya, John A. Wolf, John E. Duda, and D. Kacy Cullen. Emerging regenerative medicine

- and tissue engineering strategies for Parkinsons disease. *npj Parkinson's Disease*, 6(1):1–14, January 2020b. ISSN 2373-8057. doi: 10.1038/s41531-019-0105-5. URL <https://www.nature.com/articles/s41531-019-0105-5>. Number: 1 Publisher: Nature Publishing Group.
- Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition*. Springer Series in Statistics. Springer-Verlag, New York, 2 edition, 2009. ISBN 978-0-387-84857-0. doi: 10.1007/978-0-387-84858-7. URL <https://www.springer.com/gp/book/9780387848570>.
- M. A. Hely, J. G. Morris, R. Traficante, W. G. Reid, D. J. O'Sullivan, and P. M. Williamson. The sydney multicentre study of Parkinson's disease: progression and mortality at 10 years. *Journal of Neurology, Neurosurgery, and Psychiatry*, 67(3):300–307, September 1999. ISSN 0022-3050. doi: 10.1136/jnnp.67.3.300.
- Miguel A. Hernán, Bahi Takkouche, Francisco Caamaño-Isorna, and Juan J. Gestal-Otero. A meta-analysis of coffee drinking, cigarette smoking, and the risk of Parkinson's disease. *Annals of Neurology*, 52(3):276–284, 2002. ISSN 1531-8249. doi: 10.1002/ana.10277. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/ana.10277>. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/ana.10277>.
- Todd M. Herrington, Jennifer J. Cheng, and Emad N. Eskandar. Mechanisms of deep brain stimulation. *Journal of Neurophysiology*, 115(1):19–38, January 2016. ISSN 1522-1598. doi: 10.1152/jn.00281.2015.
- Jon P. Hiseman and Robin Fackrell. Caregiver Burden and the Nonmotor Symptoms of Parkinson's Disease. *International Review of Neurobiology*, 133:479–497, 2017. ISSN 2162-5514. doi: 10.1016/bs.irn.2017.05.035.
- S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, November 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735.
- Janet Hoenicka, Pedro J. García-Ruiz, Guillermo Ponce, Antonio Herranz, Dolores Martínez-Rubio, Estela Pérez-Santamarina, and Francesc Palau. The addiction-related gene ANKK1 in Parkinsonian patients with impulse control disorder. *Neurotoxicity Research*, 27(3):205–208, April 2015. ISSN 1476-3524. doi: 10.1007/s12640-014-9504-x.
- Bernd Holdorff. Fritz Heinrich Lewy (1885/1950). *Journal of Neurology*, 253(5):677–678, May 2006. ISSN 1432-1459. doi: 10.1007/s00415-006-0130-2. URL <https://doi.org/10.1007/s00415-006-0130-2>.
- David M. Howard, Mark J. Adams, Toni-Kim Clarke, Jonathan D. Hafferty, Jude Gibson, Masoud Shirali, Jonathan R. I. Coleman, Saskia P. Hagenaars, Joey Ward,

- Eleanor M. Wigmore, Clara Alloza, Xueyi Shen, Miruna C. Barbu, Eileen Y. Xu, Heather C. Whalley, Riccardo E. Marioni, David J. Porteous, Gail Davies, Ian J. Deary, Gibran Hemani, Klaus Berger, Henning Teismann, Rajesh Rawal, Volker Arolt, Bernhard T. Baune, Udo Dannlowski, Katharina Domschke, Chao Tian, David A. Hinds, 23andMe Research Team, Major Depressive Disorder Working Group of the Psychiatric Genomics Consortium, Maciej Trzaskowski, Enda M. Byrne, Stephan Ripke, Daniel J. Smith, Patrick F. Sullivan, Naomi R. Wray, Gerome Breen, Cathryn M. Lewis, and Andrew M. McIntosh. Genome-wide meta-analysis of depression identifies 102 independent variants and highlights the importance of the prefrontal brain regions. *Nature Neuroscience*, 22(3):343–352, 2019. ISSN 1546-1726. doi: 10.1038/s41593-018-0326-7.
- A. J. Hughes, S. E. Daniel, L. Kilford, and A. J. Lees. Accuracy of clinical diagnosis of idiopathic Parkinson’s disease: a clinico-pathological study of 100 cases. *Journal of Neurology, Neurosurgery, and Psychiatry*, 55(3):181–184, March 1992. ISSN 0022-3050. doi: 10.1136/jnnp.55.3.181.
- John D. Hunter. Matplotlib: A 2D Graphics Environment. *Computing in Science Engineering*, 9(3):90–95, May 2007. ISSN 1558-366X. doi: 10.1109/MCSE.2007.55. Conference Name: Computing in Science Engineering.
- Catherine S. Hurt, Fadi Alkufri, Richard G. Brown, David J. Burn, John V. Hindle, Sabine Landau, Kenneth C. Wilson, Michael Samuel, and PROMS-PD study group. Motor phenotypes, medication and mood: further associations with impulsive behaviours in Parkinson’s disease. *Journal of Parkinson’s Disease*, 4(2):245–254, 2014. ISSN 1877-718X. doi: 10.3233/JPD-130314.
- J. Ihle, F. Artaud, S. Bekadar, G. Mangone, S. Sambin, L. L. Mariani, H. Bertrand, O. Rascol, F. Durif, P. Derkinderen, C. Scherzer, A. Elbaz, J. C. Corvol, DIGPD study group Steering committee, Jean-Christophe Corvol, Alexis Elbaz, Marie Vidailhet, Alexis Brice, Statistical analyses, Alexis Elbaz, Fanny Artaud, Principal investigators for sites, Frédéric Bourdain, Jean-Philippe Brandel, Jean-Christophe Corvol, Pascal Derkinderen, Franck Durif, Richard Levy, Fernando Pico, Olivier Rascol, Co-investigators (alphabetical order), Anne-Marie Bonnet, Cecilia Bonnet, Christine Brefel-Courbon, Florence Cormier-Dequaire, Bertrand Degos, Bérangère Debilly, Alexis Elbaz, Monique Galitsky, David Grabli, Andreas Hartmann, Stephan Klebe, Julia Kraemmer, Lucette Lacomblez, Sara Leder, Graziella Mangone, Louise-Laure Mariani, Ana-Raquel Marques, Valérie Mesnage, Julia Muellner, Fabienne Ory-Magne, Violaine Planté-Bordeneuve, Emmanuel Roze, Melissa Tir, Marie Vidailhet, Hana You, Neuropsychologists, Eve Benchetrit, Julie Socha, Fanny Pineau, Tiphaine Vidal, Elsa Pomies, Virginie Bayet, Genetic core, Alexis Brice, Suzanne Lesage, Khadija Tahiri, Hélène Bertrand, Graziella Mangone, Sponsor activities and clini-

- cal research assistants, Alain Mallet, Coralie Villeret, Merry Mazmanian, Hakima Manseur, Mostafa Hajji, Benjamin Le Toullec, Vanessa Brochard, Monica Roy, Isabelle Rieu, Stéphane Bernard, and Antoine Faurie-Grepon. Parkinson's disease polygenic risk score is not associated with impulse control disorders: A longitudinal study. *Parkinsonism & Related Disorders*, 75:30–33, May 2020. ISSN 1873-5126. doi: 10.1016/j.parkreldis.2020.03.017.
- International League Against Epilepsy Consortium on Complex Epilepsies. Genome-wide mega-analysis identifies 16 loci and highlights diverse biological mechanisms in the common epilepsies. *Nature Communications*, 9(1):5269, 2018. ISSN 2041-1723. doi: 10.1038/s41467-018-07524-z.
- International Obsessive Compulsive Disorder Foundation Genetics Collaborative (IOCDF-GC) and OCD Collaborative Genetics Association Studies (OC GAS). Revealing the complex genetic architecture of obsessive-compulsive disorder using meta-analysis. *Molecular Psychiatry*, 23(5):1181–1188, 2018. ISSN 1476-5578. doi: 10.1038/mp.2017.154.
- A. G. Ivakhnenko and Valentin Grigorevich Lapa. *Cybernetics and Forecasting Techniques*. American Elsevier Publishing Company, 1967. Google-Books-ID: rGF-gAAAAMAAJ.
- J. Jankovic. Parkinson's disease: clinical features and diagnosis. *Journal of Neurology, Neurosurgery, and Psychiatry*, 79(4):368–376, April 2008. ISSN 1468-330X. doi: 10.1136/jnnp.2007.131045.
- Joseph Jankovic and Roger Kurlan. Tourette syndrome: evolving concepts. *Movement Disorders: Official Journal of the Movement Disorder Society*, 26(6):1149–1156, May 2011. ISSN 1531-8257. doi: 10.1002/mds.23618.
- S. Jesús, M. T. Perinán, C. Cortés, D. Buiza-Rueda, D. Macías-García, A. Adarmes, L. Muñoz-Delgado, M. Á Labrador-Espinosa, C. Tejera-Parrado, M. P. Gómez-Garre, and P. Mir. Integrating genetic and clinical data to predict impulse control disorders in Parkinson's disease. *European Journal of Neurology*, October 2020. ISSN 1468-1331. doi: 10.1111/ene.14590.
- Emma C. Johnson, Richard Border, Whitney E. Melroy-Greif, Christiaan de Leeuw, Marissa A. Ehringer, and Matthew C. Keller. No evidence that schizophrenia candidate genes are more associated with schizophrenia than non-candidate genes. *Biological psychiatry*, 82(10):702–708, November 2017. ISSN 0006-3223. doi: 10.1016/j.biopsych.2017.06.033. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5643230/>.

Juho Joutsa, Kirsti Martikainen, Tero Vahlberg, Valerie Voon, and Valtteri Kaasinen. Impulse control disorders and depression in Finnish patients with Parkinson's disease. *Parkinsonism & Related Disorders*, 18(2):155–160, February 2012. ISSN 1873-5126. doi: 10.1016/j.parkreldis.2011.09.007.

Lorraine V. Kalia and Anthony E. Lang. Parkinson's disease. *Lancet (London, England)*, 386(9996):896–912, August 2015. ISSN 1474-547X. doi: 10.1016/S0140-6736(14)61393-3.

Richard Karlsson Linnér, Pietro Biroli, Edward Kong, S. Fleur W. Meddens, Robbee Wedow, Mark Alan Fontana, Maël Lebreton, Stephen P. Tino, Abdel Abdellaoui, Anke R. Hammerschlag, Michel G. Nivard, Aysu Okbay, Cornelius A. Rietveld, Pascal N. Timshel, Maciej Trzaskowski, Ronald de Vlaming, Christian L. Zünd, Yanchun Bao, Laura Buzdugan, Ann H. Caplin, Chia-Yen Chen, Peter Eibich, Pierre Fontanillas, Juan R. Gonzalez, Peter K. Joshi, Ville Karhunen, Aaron Kleinman, Remy Z. Levin, Christina M. Lill, Gerardus A. Meddens, Gerard Muntané, Sandra Sanchez-Roige, Frank J. van Rooij, Erdogan Taskesen, Yang Wu, Futao Zhang, 23and Me Research Team, eQTLgen Consortium, International Cannabis Consortium, Social Science Genetic Association Consortium, Adam Auton, Jason D. Boardman, David W. Clark, Andrew Conlin, Conor C. Dolan, Urs Fischbacher, Patrick J. F. Groenen, Kathleen Mullan Harris, Gregor Hasler, Albert Hofman, Mohammad A. Ikram, Sonia Jain, Robert Karlsson, Ronald C. Kessler, Maarten Kooyman, James MacKillop, Minna Männikkö, Carlos Morcillo-Suarez, Matthew B. McQueen, Klaus M. Schmidt, Melissa C. Smart, Matthias Sutter, A. Roy Thurik, André G. Uitterlinden, Jon White, Harriet de Wit, Jian Yang, Lars Bertram, Dorret I. Boomsma, Tõnu Esko, Ernst Fehr, David A. Hinds, Magnus Johannesson, Meena Kumari, David Laibson, Patrik K. E. Magnusson, Michelle N. Meyer, Arcadi Navarro, Abraham A. Palmer, Tune H. Pers, Danielle Posthuma, Daniel Schunk, Murray B. Stein, Rauli Svento, Henning Tiemeier, Paul R. H. J. Timmers, Patrick Turley, Robert J. Ursano, Gert G. Wagner, James F. Wilson, Jacob Gratten, James J. Lee, David Cesarini, Daniel J. Benjamin, Philipp D. Koellinger, and Jonathan P. Beauchamp. Genome-wide association analyses of risk tolerance and risky behaviors in over 1 million individuals identify hundreds of loci and shared genetic influences. *Nature Genetics*, 51(2):245–257, 2019. ISSN 1546-1718. doi: 10.1038/s41588-018-0309-3.

Gülay Kenangil, Sibel Ozekmekçi, Melis Sohtaoglu, and Ethem Erginöz. Compulsive behaviors in patients with Parkinson's disease. *The Neurologist*, 16(3):192–195, May 2010. ISSN 2331-2637. doi: 10.1097/NRL.0b013e31819f952b.

Hee J. Kim, Beom S. Jeon, and Peter Jenner. Hallmarks of Treatment Aspects: Parkinson's Disease Throughout Centuries Including l-Dopa. *International Review of Neurobiology*, 132:295–343, 2017. ISSN 2162-5514. doi: 10.1016/bs.irn.2017.01.006.

- Atesh Koul, Cristina Becchio, and Andrea Cavallo. Cross-Validation Approaches for Replicability in Psychology. *Frontiers in Psychology*, 9, 2018. ISSN 1664-1078. doi: 10.3389/fpsyg.2018.01117. URL <https://www.frontiersin.org/articles/10.3389/fpsyg.2018.01117/full>. Publisher: Frontiers.
- Julia Kraemmer, Kara Smith, Daniel Weintraub, Vincent Guillemot, Mike A. Nalls, Florence Cormier-Dequaire, Ivan Moszer, Alexis Brice, Andrew B. Singleton, and Jean-Christophe Corvol. Clinical-genetic model predicts incident impulse control disorders in Parkinson’s disease. *Journal of Neurology, Neurosurgery, and Psychiatry*, 87(10):1106–1111, October 2016. ISSN 1468-330X. doi: 10.1136/jnnp-2015-312848.
- Soumya Krishnamoorthy, Roopa Rajan, Moinak Banerjee, Hardeep Kumar, Gangadhara Sarma, Syam Krishnan, Sankara Sarma, and Asha Kishore. Dopamine D3 receptor Ser9Gly variant is associated with impulse control disorders in Parkinson’s disease patients. *Parkinsonism & Related Disorders*, 30:13–17, 2016. ISSN 1873-5126. doi: 10.1016/j.parkreldis.2016.06.005.
- Jeppe Theiss Kristensen and Paolo Burelli. Combining Sequential and Aggregated Data for Churn Prediction in Casual Freemium Games. In *2019 IEEE Conference on Games (CoG)*, pages 1–8, London, United Kingdom, August 2019. IEEE. ISBN 978-1-72811-884-0. doi: 10.1109/CIG.2019.8848106. URL <https://ieeexplore.ieee.org/document/8848106/>.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet Classification with Deep Convolutional Neural Networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012. URL <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>.
- Jee-Young Lee, Eun Kyung Lee, Sung Sup Park, Ji-Yeon Lim, Hee Jin Kim, Ji Sun Kim, and Beom S. Jeon. Association of DRD3 and GRIN2B with impulse control and related behaviors in Parkinson’s disease. *Movement Disorders: Official Journal of the Movement Disorder Society*, 24(12):1803–1810, September 2009. ISSN 1531-8257. doi: 10.1002/mds.22678.
- Jee-Young Lee, Jong-Min Kim, Jae Woo Kim, Jinwhan Cho, Won Yong Lee, Han-Joon Kim, and Beom S. Jeon. Association between the dose of dopaminergic medication and the behavioral disturbances in Parkinson disease. *Parkinsonism & Related Disorders*, 16(3):202–207, March 2010. ISSN 1873-5126. doi: 10.1016/j.parkreldis.2009.12.002.

- Jee-Young Lee, Beom S. Jeon, Han-Joon Kim, and Sung-Sup Park. Genetic variant of HTR2A associates with risk of impulse control and repetitive behaviors in Parkinson's disease. *Parkinsonism & Related Disorders*, 18(1):76–78, January 2012. ISSN 1873-5126. doi: 10.1016/j.parkreldis.2011.08.009.
- Andrew J. Lees. Unresolved issues relating to the shaking palsy on the celebration of James Parkinson's 250th birthday. *Movement Disorders: Official Journal of the Movement Disorder Society*, 22 Suppl 17:S327–334, September 2007. ISSN 0885-3185. doi: 10.1002/mds.21684.
- Anna Leontjeva and Ilya Kuzovkin. Combining Static and Dynamic Features for Multivariate Sequence Classification. *2016 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, pages 21–30, October 2016. doi: 10.1109/DSAA.2016.10. URL <http://arxiv.org/abs/1712.08160>. arXiv: 1712.08160.
- Eric W. Leppink, Katherine Lust, and Jon E. Grant. Depression in university students: associations with impulse control disorders. *International Journal of Psychiatry in Clinical Practice*, 20(3):146–150, September 2016a. ISSN 1471-1788. doi: 10.1080/13651501.2016.1197272.
- Eric W. Leppink, Brian L. Odlaug, Katherine Lust, Gary Christenson, and Jon E. Grant. The Young and the Stressed: Stress, Impulse Control, and Health in College Students. *The Journal of Nervous and Mental Disease*, 204(12):931–938, December 2016b. ISSN 1539-736X. doi: 10.1097/NMD.0000000000000586.
- Iracema Leroi, Michelle Andrews, Kathryn McDonald, Vijay Harbishettar, Rebecca Elliott, E. Jane Byrne, and Alistair Burns. Apathy and impulse control disorders in Parkinson's disease: a direct comparison. *Parkinsonism & Related Disorders*, 18(2):198–203, February 2012. ISSN 1873-5126. doi: 10.1016/j.parkreldis.2011.10.005.
- Smaranda Leu-Semenescu, Elias Karroum, Agnès Brion, Eric Konofal, and Isabelle Arnulf. Dopamine dysregulation syndrome in a patient with restless legs syndrome. *Sleep Medicine*, 10(4):494–496, April 2009. ISSN 1389-9457. doi: 10.1016/j.sleep.2008.12.010.
- Irene Litvan, Dag Aarsland, Charles H. Adler, Jennifer G. Goldman, Jaime Kulisevsky, Brit Mollenhauer, Maria C. Rodriguez-Oroz, Alexander I. Tröster, and Daniel Weintraub. MDS Task Force on mild cognitive impairment in Parkinson's disease: critical review of PD-MCI. *Movement Disorders: Official Journal of the Movement Disorder Society*, 26(10):1814–1824, August 2011. ISSN 1531-8257. doi: 10.1002/mds.23823.
- Ganqiang Liu, Joseph J. Locascio, Jean-Christophe Corvol, Brendon Boot, Zhixiang Liao, Kara Page, Daly Franco, Kyle Burke, Iris E. Jansen, Ana Trisini-Lipsanopoulos, Sophie Winder-Rhodes, Caroline M. Tanner, Anthony E. Lang, Shirley Eberly,

- Alexis Elbaz, Alexis Brice, Graziella Mangone, Bernard Ravina, Ira Shoulson, Florence Cormier-Dequaire, Peter Heutink, Jacobus J. van Hilten, Roger A. Barker, Caroline H. Williams-Gray, Johan Marinus, Clemens R. Scherzer, HBS, CamPaIGN, PICNICS, PROPARK, PSG, DIGPD, and PDBP. Prediction of cognition in Parkinson's disease with a clinical-genetic score: a longitudinal analysis of nine cohorts. *The Lancet. Neurology*, 16(8):620–629, 2017. ISSN 1474-4465. doi: 10.1016/S1474-4422(17)30122-9.
- Mengzhen Liu, Yu Jiang, Robbee Wedow, Yue Li, David M. Brazeal, Fang Chen, Gargi Datta, Jose Davila-Velderrain, Daniel McGuire, Chao Tian, Xiaowei Zhan, 23andMe Research Team, HUNT All-In Psychiatry, Hélène Choquet, Anna R. Docherty, Jessica D. Faul, Johanna R. Foerster, Lars G. Fritsche, Maiken Elvestad Gabrielsen, Scott D. Gordon, Jeffrey Haessler, Jouke-Jan Hottenga, Hongyan Huang, Seon-Kyeong Jang, Philip R. Jansen, Yueh Ling, Reedik Mägi, Nana Matoba, George McMahon, Antonella Mulas, Valeria Orrù, Teemu Palviainen, Anita Pandit, Gunnar W. Reginsson, Anne Heidi Skogholt, Jennifer A. Smith, Amy E. Taylor, Constance Turman, Gonneke Willemsen, Hannah Young, Kendra A. Young, Gregory J. M. Zalc, Wei Zhao, Wei Zhou, Gyda Bjornsdottir, Jason D. Boardman, Michael Boehnke, Dorret I. Boomsma, Chu Chen, Francesco Cucca, Gareth E. Davies, Charles B. Eaton, Marissa A. Ehringer, Tõnu Esko, Edoardo Fiorillo, Nathan A. Gillespie, Daniel F. Gudbjartsson, Toomas Haller, Kathleen Mullan Harris, Andrew C. Heath, John K. Hewitt, Ian B. Hickie, John E. Hokanson, Christian J. Hopfer, David J. Hunter, William G. Iacono, Eric O. Johnson, Yoichiro Kamatani, Sharon L. R. Kardia, Matthew C. Keller, Manolis Kellis, Charles Kooperberg, Peter Kraft, Kenneth S. Krauter, Markku Laakso, Penelope A. Lind, Anu Loukola, Sharon M. Lutz, Pamela A. F. Madden, Nicholas G. Martin, Matt McGue, Matthew B. McQueen, Sarah E. Medland, Andres Metspalu, Karen L. Mohlke, Jonas B. Nielsen, Yukinori Okada, Ulrike Peters, Tinca J. C. Polderman, Danielle Posthuma, Alexander P. Reiner, John P. Rice, Eric Rimm, Richard J. Rose, Valgerdur Runarsdottir, Michael C. Stallings, Alena Stanáková, Hreinn Stefansson, Khanh K. Thai, Hilary A. Tindle, Thorarinn Tyrfinngsson, Tamara L. Wall, David R. Weir, Constance Weisner, John B. Whitfield, Bendik Slagsvold Winsvold, Jie Yin, Luisa Zuccolo, Laura J. Bierut, Kristian Hveem, James J. Lee, Marcus R. Munafò, Nancy L. Saccone, Cristen J. Willer, Marilyn C. Cornelis, Sean P. David, David A. Hinds, Eric Jorgenson, Jaakko Kaprio, Jerry A. Stitzel, Kari Stefansson, Thorgeir E. Thorgeirsson, Gonçalo Abecasis, Dajiang J. Liu, and Scott Vrieze. Association studies of up to 1.2 million individuals yield new insights into the genetic etiology of tobacco and alcohol use. *Nature Genetics*, 51(2): 237–244, 2019. ISSN 1546-1718. doi: 10.1038/s41588-018-0307-5.
- Min-Tzu Lo, David A. Hinds, Joyce Y. Tung, Carol Franz, Chun-Chieh Fan, Yunpeng Wang, Olav B. Smeland, Andrew Schork, Dominic Holland, Karolina Kauppi,

- Nilotpal Sanyal, Valentina Escott-Price, Daniel J. Smith, Michael O'Donovan, Hreinn Stefansson, Gyda Bjornsdottir, Thorgeir E. Thorgeirsson, Kari Stefansson, Linda K. McEvoy, Anders M. Dale, Ole A. Andreassen, and Chi-Hua Chen. Genome-wide analyses for personality traits identify six genomic loci and show correlations with psychiatric disorders. *Nature Genetics*, 49(1):152–156, 2017. ISSN 1546-1718. doi: 10.1038/ng.3736.
- Michelle Luciano, Saskia P. Hagenaars, Gail Davies, W. David Hill, Toni-Kim Clarke, Masoud Shirali, Sarah E. Harris, Riccardo E. Marioni, David C. Liewald, Chloe Fawns-Ritchie, Mark J. Adams, David M. Howard, Cathryn M. Lewis, Catharine R. Gale, Andrew M. McIntosh, and Ian J. Deary. Association analysis in over 329,000 individuals identifies 116 independent variants influencing neuroticism. *Nature Genetics*, 50(1):6–11, 2018. ISSN 1546-1718. doi: 10.1038/s41588-017-0013-8.
- M. R. Luquin, O. Scipioni, J. Vaamonde, O. Gershanik, and J. A. Obeso. Levodopa-induced dyskinesias in Parkinson's disease: clinical and pharmacological classification. *Movement Disorders: Official Journal of the Movement Disorder Society*, 7(2):117–124, 1992. ISSN 0885-3185. doi: 10.1002/mds.870070204.
- Eugenia Mamikonyan, Andrew D. Siderowf, John E. Duda, Marc N. Potenza, Stacy Horn, Matthew B. Stern, and Daniel Weintraub. Long-term follow-up of impulse control disorders in Parkinson's disease. *Movement Disorders*, 23(1):75–80, 2008. ISSN 1531-8257. doi: 10.1002/mds.21770. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/mds.21770>. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/mds.21770>.
- Kenneth Marek, Danna Jennings, Shirley Lasch, Andrew Siderowf, Caroline Tanner, Tanya Simuni, Chris Coffey, Karl Kieburtz, Emily Flagg, Sohini Chowdhury, Werner Poewe, Brit Mollenhauer, Paracelsus-Elena Klinik, Todd Sherer, Mark Frasier, Claire Meunier, Alice Rudolph, Cindy Casaceli, John Seibyl, Susan Mendick, Norbert Schuff, Ying Zhang, Arthur Toga, Karen Crawford, Alison Ansbach, Pasquale De Blasio, Michele Piovella, John Trojanowski, Les Shaw, Andrew Singleton, Keith Hawkins, Jamie Eberling, Deborah Brooks, David Russell, Laura Leary, Stewart Factor, Barbara Sommerfeld, Penelope Hogarth, Emily Pighetti, Karen Williams, David Standaert, Stephanie Guthrie, Robert Hauser, Holly Delgado, Joseph Jankovic, Christine Hunter, Matthew Stern, Baochan Tran, Jim Leverenz, Marne Baca, Sam Frank, Cathi-Ann Thomas, Irene Richard, Cheryl Deeley, Linda Rees, Fabienne Sprenger, Elisabeth Lang, Holly Shill, Sanja Obradov, Hubert Fernandez, Adrienna Winters, Daniela Berg, Katharina Gauss, Douglas Galasko, Deborah Fontaine, Zoltan Mari, Melissa Gerstenhaber, David Brooks, Sophie Malloy, Paolo Barone, Katia Longo, Tom Comery, Bernard Ravina, Igor Grachev, Kim Gallagher, Michelle Collins, Katherine L. Widnell, Suzanne Ostrowizki, Paulo Fontoura, Tony Ho, Johan Luth-

- man, Marcel van der Brug, Alastair D. Reith, and Peggy Taylor. The Parkinson Progression Marker Initiative (PPMI). *Progress in Neurobiology*, 95(4):629–635, December 2011. ISSN 0301-0082. doi: 10.1016/j.pneurobio.2011.09.005. URL <http://www.sciencedirect.com/science/article/pii/S0301008211001651>.
- Ana Marques, Michela Figorilli, Bruno Pereira, Philippe Derost, Berengere Debilly, Patricia Beudin, Tiphaine Vidal, Franck Durif, and Maria Livia Fantini. Impulse control disorders in Parkinson’s disease patients with RLS: a cross sectional-study. *Sleep Medicine*, 48:148–154, 2018. ISSN 1878-5506. doi: 10.1016/j.sleep.2018.02.004.
- Ana Marques, Tiphaine Vidal, Bruno Pereira, Eve Benchetrit, Julie Socha, Fanny Pineau, Alexis Elbaz, Fanny Artaud, Graziella Mangone, Hana You, Florence Cormier, Monique Galitstky, Elsa Pomies, Olivier Rascol, Pascal Derkinderen, Daniel Weintraub, Jean Christophe Corvol, Franck Durif, and DIGPD study group. French validation of the questionnaire for Impulsive-Compulsive Disorders in Parkinson’s Disease-Rating Scale (QUIP-RS). *Parkinsonism & Related Disorders*, 63:117–123, 2019. ISSN 1873-5126. doi: 10.1016/j.parkreldis.2019.02.026.
- Alicia R. Martin, Christopher R. Gignoux, Raymond K. Walters, Genevieve L. Wojcik, Benjamin M. Neale, Simon Gravel, Mark J. Daly, Carlos D. Bustamante, and Eimear E. Kenny. Human Demographic History Impacts Genetic Risk Prediction across Diverse Populations. *The American Journal of Human Genetics*, 100(4):635–649, April 2017. ISSN 0002-9297. doi: 10.1016/j.ajhg.2017.03.004. URL <http://www.sciencedirect.com/science/article/pii/S0002929717301076>.
- Llew Mason, Jonathan Baxter, Peter L. Bartlett, and Marcus R. Frean. Boosting Algorithms as Gradient Descent. In S. A. Solla, T. K. Leen, and K. Müller, editors, *Advances in Neural Information Processing Systems 12*, pages 512–518. MIT Press, 2000. URL <http://papers.nips.cc/paper/1766-boosting-algorithms-as-gradient-descent.pdf>.
- Nozomu Matsuda, Shunsuke Kobayashi, and Yoshikazu Ugawa. [Devotion to painting in a Parkinson’s disease patient]. *Rinsho Shinkeigaku = Clinical Neurology*, 58(12):756–760, December 2018. ISSN 1882-0654. doi: 10.5692/clinicalneurol.cn-001182.
- L. Mazzella, M. D. Yahr, L. Marinelli, N. Huang, E. Moshier, and A. Di Rocco. Dyskinesias predict the onset of motor response fluctuations in patients with Parkinson’s disease on l-dopa monotherapy. *Parkinsonism & Related Disorders*, 11(3):151–155, May 2005. ISSN 1353-8020, 1873-5126. doi: 10.1016/j.parkreldis.2004.10.002. URL [https://www.prd-journal.com/article/S1353-8020\(04\)00158-0/abstract](https://www.prd-journal.com/article/S1353-8020(04)00158-0/abstract). Publisher: Elsevier.
- Shane McCarthy, Sayantan Das, Warren Kretschmar, Olivier Delaneau, Andrew R. Wood, Alexander Teumer, Hyun Min Kang, Christian Fuchsberger, Petr Danecek,

- Kevin Sharp, Yang Luo, Carlo Sidore, Alan Kwong, Nicholas Timpson, Seppo Koskinen, Scott Vrieze, Laura J. Scott, He Zhang, Anubha Mahajan, Jan Veldink, Ulrike Peters, Carlos Pato, Cornelia M. van Duijn, Christopher E. Gillies, Ilaria Gandin, Massimo Mezzavilla, Arthur Gilly, Massimiliano Cocca, Michela Traglia, Andrea Angius, Jeffrey C. Barrett, Dorrett Boomsma, Kari Branham, Gerome Breen, Chad M. Brummett, Fabio Busonero, Harry Campbell, Andrew Chan, Sai Chen, Emily Chew, Francis S. Collins, Laura J. Corbin, George Davey Smith, George Dedoussis, Marcus Dorr, Aliko-Eleni Farmaki, Luigi Ferrucci, Lukas Forer, Ross M. Fraser, Stacey Gabriel, Shawn Levy, Leif Groop, Tabitha Harrison, Andrew Hattersley, Oddgeir L. Holmen, Kristian Hveem, Matthias Kretzler, James C. Lee, Matt McGue, Thomas Meitinger, David Melzer, Josine L. Min, Karen L. Mohlke, John B. Vincent, Matthias Nauck, Deborah Nickerson, Aarno Palotie, Michele Pato, Nicola Pirastu, Melvin McInnis, J. Brent Richards, Cinzia Sala, Veikko Salomaa, David Schlessinger, Sebastian Schoenherr, P. Eline Slagboom, Kerrin Small, Timothy Spector, Dwight Stambolian, Marcus Tuke, Jaakko Tuomilehto, Leonard H. Van den Berg, Wouter Van Rheenen, Uwe Volker, Cisca Wijmenga, Daniela Toniolo, Eleftheria Zeggini, Paolo Gasparini, Matthew G. Sampson, James F. Wilson, Timothy Frayling, Paul I. W. de Bakker, Morris A. Swertz, Steven McCarroll, Charles Kooperberg, Annelot Dekker, David Altshuler, Cristen Willer, William Iacono, Samuli Ripatti, Nicole Soranzo, Klaudia Walter, Anand Swaroop, Francesco Cucca, Carl A. Anderson, Richard M. Myers, Michael Boehnke, Mark I. McCarthy, Richard Durbin, and Haplotype Reference Consortium. A reference panel of 64,976 haplotypes for genotype imputation. *Nature Genetics*, 48(10):1279–1283, 2016. ISSN 1546-1718. doi: 10.1038/ng.3643.
- S. L. McElroy, H. G. Pope, P. E. Keck, J. I. Hudson, K. A. Phillips, and S. M. Strakowski. Are impulse-control disorders related to bipolar disorder? *Comprehensive Psychiatry*, 37(4):229–240, August 1996. ISSN 0010-440X. doi: 10.1016/s0010-440x(96)90001-2.
- Leland McInnes, John Healy, and James Melville. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *arXiv:1802.03426 [cs, stat]*, December 2018. URL <http://arxiv.org/abs/1802.03426>. arXiv: 1802.03426.
- Wes McKinney. Data Structures for Statistical Computing in Python. *Proceedings of the 9th Python in Science Conference*, pages 56–61, 2010. doi: 10.25080/Majora-92bf1922-00a. URL <http://conference.scipy.org/proceedings/scipy2010/mckinney.html>. Conference Name: Proceedings of the 9th Python in Science Conference.
- T. H. Mertsalmi, V. T. E. Aho, P. a. B. Pereira, L. Paulin, E. Pekkonen, P. Auvinen, and F. Scheperjans. More than constipation - bowel symptoms in Parkinson’s disease

- and their connection to gut microbiota. *European Journal of Neurology*, 24(11):1375–1383, 2017. ISSN 1468-1331. doi: 10.1111/ene.13398.
- Raymond G. Miltenberger, Jennifer Redlin, Ross Crosby, Marcella Stickney, Jim Mitchell, Stephen Wonderlich, Ronald Faber, and Joshua Smyth. Direct and retrospective assessment of factors contributing to compulsive buying. *Journal of Behavior Therapy and Experimental Psychiatry*, 34(1):1–9, March 2003. ISSN 0005-7916. doi: 10.1016/s0005-7916(03)00002-8.
- Takayasu Mishima, Shinsuke Fujioka, Ryoichi Kurisaki, Shozaburo Yanamoto, Masa-aki Higuchi, Jun Tsugawa, Jiro Fukae, Ryuji Neshige, and Yoshio Tsuboi. Impulse control disorders and punding in Perry syndrome. *Parkinsonism & Related Disorders*, 21(11):1381–1382, November 2015. ISSN 1873-5126. doi: 10.1016/j.parkreldis.2015.09.037.
- Pooya Mobadersany, Safoora Yousefi, Mohamed Amgad, David A. Gutman, Jill S. Barnholtz-Sloan, José E. Velázquez Vega, Daniel J. Brat, and Lee A. D. Cooper. Predicting cancer outcomes from histology and genomics using convolutional networks. *Proceedings of the National Academy of Sciences*, 115(13):E2970–E2979, March 2018. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.1717139115. URL <https://www.pnas.org/content/115/13/E2970>. ZSCC: 0000143 Publisher: National Academy of Sciences Section: PNAS Plus.
- Thomas J. Moore, Joseph Glenmullen, and Donald R. Mattison. Reports of pathological gambling, hypersexuality, and compulsive shopping associated with dopamine receptor agonist drugs. *JAMA internal medicine*, 174(12):1930–1933, December 2014. ISSN 2168-6114. doi: 10.1001/jamainternmed.2014.5262.
- Ahmed A. Moustafa, Srinivasa Chakravarthy, Joseph R. Phillips, Ankur Gupta, Szabolcs Keri, Bertalan Polner, Michael J. Frank, and Marjan Jahanshahi. Motor symptoms in Parkinson’s disease: A unified framework. *Neuroscience and Biobehavioral Reviews*, 68:727–740, September 2016. ISSN 1873-7528. doi: 10.1016/j.neubiorev.2016.07.010.
- M. D. Muentner and G. M. Tyce. L-dopa therapy of Parkinson’s disease: plasma L-dopa concentration, therapeutic response, and side effects. *Mayo Clinic Proceedings*, 46(4):231–239, April 1971. ISSN 0025-6196.
- Kevin P. Murphy. *Machine Learning: A Probabilistic Perspective*. The MIT Press, 2012. ISBN 978-0-262-01802-9.
- Mike A. Nalls, Jose Bras, Dena G. Hernandez, Margaux F. Keller, Elisa Majounie, Alan E. Renton, Mohamad Saad, Iris Jansen, Rita Guerreiro, Steven Lubbe, Vincent Plagnol, J. Raphael Gibbs, Claudia Schulte, Nathan Pankratz, Margaret Sutherland, Lars Bertram, Christina M. Lill, Anita L. DeStefano, Tatiana Faroud, Nicholas

- Eriksson, Joyce Y. Tung, Connor Edsall, Noah Nichols, Janet Brooks, Sampath Arepalli, Hannah Pliner, Chris Letson, Peter Heutink, Maria Martinez, Thomas Gasser, Bryan J. Traynor, Nick Wood, John Hardy, and Andrew B. Singleton. NeuroX, a fast and efficient genotyping platform for investigation of neurodegenerative diseases. *Neurobiology of aging*, 36(3):1605.e7–1605.12, March 2015. ISSN 0197-4580. doi: 10.1016/j.neurobiolaging.2014.07.028. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4317375/>.
- Mike A. Nalls, Cornelis Blauwendraat, Costanza L. Vallerga, Karl Heilbron, Sara Bandres-Ciga, Diana Chang, Manuela Tan, Demis A. Kia, Alastair J. Noyce, Angli Xue, Jose Bras, Emily Young, Rainer von Coelln, Javier Simón-Sánchez, Claudia Schulte, Manu Sharma, Lynne Krohn, Lasse Pihlstrøm, Ari Siitonen, Hirotaka Iwaki, Hampton Leonard, Faraz Faghri, J. Raphael Gibbs, Dena G. Hernandez, Sonja W. Scholz, Juan A. Botia, Maria Martinez, Jean-Christophe Corvol, Suzanne Lesage, Joseph Jankovic, Lisa M. Shulman, Margaret Sutherland, Pentti Tienari, Kari Majamaa, Mathias Toft, Ole A. Andreassen, Tushar Bangale, Alexis Brice, Jian Yang, Ziv Gan-Or, Thomas Gasser, Peter Heutink, Joshua M. Shulman, Nicholas W. Wood, David A. Hinds, John A. Hardy, Huw R. Morris, Jacob Gratten, Peter M. Visscher, Robert R. Graham, Andrew B. Singleton, 23andMe Research Team, System Genomics of Parkinson’s Disease Consortium, and International Parkinson’s Disease Genomics Consortium. Identification of novel risk loci, causal insights, and heritable risk for Parkinson’s disease: a meta-analysis of genome-wide association studies. *The Lancet. Neurology*, 18(12):1091–1102, 2019. ISSN 1474-4465. doi: 10.1016/S1474-4422(19)30320-5.
- National Clinical Guideline for Diagnosis and Management in Primary and Secondary Care. Symptomatic pharmacological therapy in Parkinsons disease. In *Parkinson’s Disease*. Royal College of Physicians (UK), 2006.
- Neale lab. UK Biobank, 2018. URL <http://www.nealelab.is/uk-biobank>.
- Melissa J. Nirenberg and Cheryl Waters. Compulsive eating and weight gain related to dopamine agonist use. *Movement Disorders*, 21(4): 524–529, 2006. ISSN 1531-8257. doi: 10.1002/mds.20757. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/mds.20757>. __eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/mds.20757>.
- Robert L. Nussbaum and Christopher E. Ellis. Alzheimer’s disease and Parkinson’s disease. *The New England Journal of Medicine*, 348(14):1356–1364, April 2003. ISSN 1533-4406. doi: 10.1056/NEJM2003ra020003.
- J. G. Nutt. Levodopa-induced dyskinesia: review, observations, and speculations. *Neurology*, 40(2):340–345, February 1990. ISSN 0028-3878. doi: 10.1212/wnl.40.2.340.

- J. G. Nutt and N. H. Holford. The response to levodopa in Parkinson's disease: imposing pharmacological law and order. *Annals of Neurology*, 39(5):561–573, May 1996. ISSN 0364-5134. doi: 10.1002/ana.410390504.
- Jose A. Obeso, Maria Cruz Rodríguez-Oroz, Beatriz Benitez-Temino, Francisco J. Blesa, Jorge Guridi, Concepción Marin, and Manuel Rodriguez. Functional organization of the basal ganglia: therapeutic implications for Parkinson's disease. *Movement Disorders: Official Journal of the Movement Disorder Society*, 23 Suppl 3:S548–559, 2008. ISSN 1531-8257. doi: 10.1002/mds.22062.
- Brian L. Odlaug and Jon E. Grant. Impulse-control disorders in a college sample: results from the self-administered Minnesota Impulse Disorders Interview (MIDI). *Primary Care Companion to the Journal of Clinical Psychiatry*, 12(2), 2010. ISSN 1555-211X. doi: 10.4088/PCC.09m00842whi.
- S. Ogasahara, Y. Nishikawa, M. Takahashi, K. Wada, Y. Nakamura, S. Yorifuji, and S. Tarui. Dopamine metabolism in the central nervous system after discontinuation of L-dopa therapy in patients with Parkinson disease. *Journal of the Neurological Sciences*, 66(2-3):151–163, December 1984. ISSN 0022-510X. doi: 10.1016/0022-510x(84)90003-0.
- Jacqueline Olley, Alex Blaszczynski, and Simon Lewis. Dopaminergic Medication in Parkinson's Disease and Problem Gambling. *Journal of Gambling Studies*, 31(3): 1085–1106, September 2015. ISSN 1573-3602. doi: 10.1007/s10899-014-9503-0.
- Fabienne Ory-Magne, Jean-Christophe Corvol, Jean-Philippe Azulay, Anne-Marie Bonnet, Christine Brefel-Courbon, Philippe Damier, Estelle Dellapina, Alain Destée, Franck Durif, Monique Galitzky, Thibaud Lebouvier, Wassilios Meissner, Claire Thalamas, François Tison, Alexandrine Salis, Agnès Sommet, François Viallet, Marie Vidailhet, Olivier Rascol, and NS-Park CIC Network. Withdrawing amantadine in dyskinetic patients with Parkinson disease: the AMANDYSK trial. *Neurology*, 82(4): 300–307, January 2014. ISSN 1526-632X. doi: 10.1212/WNL.0000000000000050.
- Sean S. O'Sullivan, Clare M. Loane, Andrew D. Lawrence, Andrew H. Evans, Paola Piccini, and Andrew J. Lees. Sleep disturbance and impulsive-compulsive behaviours in Parkinson's disease. *Journal of Neurology, Neurosurgery, and Psychiatry*, 82(6): 620–622, June 2011. ISSN 1468-330X. doi: 10.1136/jnnp.2009.186874.
- T. Otowa, K. Hek, M. Lee, E. M. Byrne, S. S. Mirza, M. G. Nivard, T. Bigdeli, S. H. Aggen, D. Adkins, A. Wolen, A. Fanous, M. C. Keller, E. Castelao, Z. Kutalik, S. Van der Auwera, G. Homuth, M. Nauck, A. Teumer, Y. Milaneschi, J.-J. Hottenga, N. Direk, A. Hofman, A. Uitterlinden, C. L. Mulder, A. K. Henders, S. E. Medland, S. Gordon, A. C. Heath, P. a. F. Madden, M. L. Pergadia, P. J. van der Most, I. M.

- Nolte, F. V. A. van Oort, C. A. Hartman, A. J. Oldehinkel, M. Preisig, H. J. Grabe, C. M. Middeldorp, B. W. J. H. Penninx, D. Boomsma, N. G. Martin, G. Montgomery, B. S. Maher, E. J. van den Oord, N. R. Wray, H. Tiemeier, and J. M. Hettema. Meta-analysis of genome-wide association studies of anxiety disorders. *Molecular Psychiatry*, 21(10):1391–1399, 2016. ISSN 1476-5578. doi: 10.1038/mp.2015.197.
- James Parkinson. An Essay on the Shaking Palsy. *J Neuropsychiatry Clin Neurosci*, page 14, 2002.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d\textquotesingle Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8026–8037. Curran Associates, Inc., 2019. URL <http://papers.nips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duchesnay. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12(85):2825–2830, 2011. ISSN 1533-7928. URL <http://jmlr.org/papers/v12/pedregosa11a.html>.
- Michele Poletti, Chiara Logi, Claudio Lucetti, Paolo Del Dotto, Filippo Baldacci, Andrea Vergallo, Martina Ulivi, Simone Del Sarto, Giuseppe Rossi, Roberto Ceravolo, and Ubaldo Bonuccelli. A single-center, cross-sectional prevalence study of impulse control disorders in Parkinson disease: association with dopaminergic drugs. *Journal of Clinical Psychopharmacology*, 33(5):691–694, October 2013. ISSN 1533-712X. doi: 10.1097/JCP.0b013e3182979830.
- Francesco E. Pontieri, Francesca Assogna, Clelia Pellicano, Claudia Cacciari, Sara Pannunzi, Annalucia Morrone, Emanuela Danese, Carlo Caltagirone, and Gianfranco Spalletta. Sociodemographic, neuropsychiatric and cognitive characteristics of pathological gambling and impulse control disorders NOS in Parkinson’s disease. *European Neuropsychopharmacology: The Journal of the European College of Neuropsychopharmacology*, 25(1):69–76, January 2015. ISSN 1873-7862. doi: 10.1016/j.euroneuro.2014.11.006.

- Gregory Pontone, James R. Williams, Susan Spear Bassett, and Laur Marsh. Clinical features associated with impulse control disorders in Parkinson disease. *Neurology*, 67(7):1258–1261, October 2006. ISSN 1526-632X. doi: 10.1212/01.wnl.0000238401.76928.45.
- Ronald B. Postuma, Daniela Berg, Matthew Stern, Werner Poewe, C. Warren Olanow, Wolfgang Oertel, José Obeso, Kenneth Marek, Irene Litvan, Anthony E. Lang, Glenda Halliday, Christopher G. Goetz, Thomas Gasser, Bruno Dubois, Piu Chan, Bastiaan R. Bloem, Charles H. Adler, and Günther Deuschl. MDS clinical diagnostic criteria for Parkinson’s disease. *Movement Disorders: Official Journal of the Movement Disorder Society*, 30(12):1591–1601, October 2015. ISSN 1531-8257. doi: 10.1002/mds.26424.
- Robert A. Power, Stacy Steinberg, Gyda Bjornsdottir, Cornelius A. Rietveld, Abdel Abdellaoui, Michel M. Nivard, Magnus Johannesson, Tessel E. Galesloot, Jouke J. Hottenga, Gonneke Willemsen, David Cesarini, Daniel J. Benjamin, Patrik K. E. Magnusson, Fredrik Ullén, Henning Tiemeier, Albert Hofman, Frank J. A. van Rooij, G. Bragi Walters, Engilbert Sigurdsson, Thorgeir E. Thorgeirsson, Andres Ingason, Agnar Helgason, Augustine Kong, Lambertus A. Kiemeny, Philipp Koellinger, Dorret I. Boomsma, Daniel Gudbjartsson, Hreinn Stefansson, and Kari Stefansson. Polygenic risk scores for schizophrenia and bipolar disorder predict creativity. *Nature Neuroscience*, 18(7):953–955, July 2015. ISSN 1546-1726. doi: 10.1038/nn.4040. URL <https://www.nature.com/articles/nn.4040/>. Number: 7 Publisher: Nature Publishing Group.
- Catharina Claudia Probst, Lina Marie Winter, Bettina Möller, Heinz Weber, Daniel Weintraub, Karsten Witt, Günther Deuschl, Regina Katzenschlager, and Thilo van Eimeren. Validation of the questionnaire for impulsive-compulsive disorders in Parkinson’s disease (QUIP) and the QUIP-rating scale in a German speaking sample. *Journal of Neurology*, 261(5):936–942, May 2014. ISSN 1432-1459. doi: 10.1007/s00415-014-7299-6.
- Sara L. Pulit, Charli Stoneman, Andrew P. Morris, Andrew R. Wood, Craig A. Glastonbury, Jessica Tyrrell, Loïc Yengo, Teresa Ferreira, Eirini Marouli, Yingjie Ji, Jian Yang, Samuel Jones, Robin Beaumont, Damien C. Croteau-Chonka, Thomas W. Winkler, GIANT Consortium, Andrew T. Hattersley, Ruth J. F. Loos, Joel N. Hirschhorn, Peter M. Visscher, Timothy M. Frayling, Hanieh Yaghootkar, and Cecilia M. Lindgren. Meta-analysis of genome-wide association studies for body fat distribution in 694,649 individuals of European ancestry. *Human Molecular Genetics*, 28(1):166–174, 2019. ISSN 1460-2083. doi: 10.1093/hmg/ddy327.
- A. Punjabi, A. Martersteck, Y. Wang, T.B. Parrish, and A.K. Katsaggelos. Neuroimag-

- ing modality fusion in Alzheimer’s classification using convolutional neural networks. *PLoS ONE*, 14(12), 2019. doi: 10.1371/journal.pone.0225759. ZSCC: 0000002.
- Molla Hafizur Rahman, Shuhan Yuan, Charles Xie, and Zhenghui Sha. Predicting human design decisions with deep recurrent neural network combining static and dynamic data. *Design Science*, 6, 2020. ISSN 2053-4701. doi: 10.1017/dsj.2020.12. URL <https://www.cambridge.org/core/journals/design-science/article/predicting-human-design-decisions-with-deep-recurrent-neural-network-combining-static-and-dynamic-data/097456E3CE09F11435F535B507AE9B8B>. Publisher: Cambridge University Press.
- Carolina Candelaria Ramírez Gómez, Marcos Serrano Dueñas, Oscar Bernal, Natalia Araoz, Michel Sáenz Farret, Victoria Aldinio, Verónica Montilla, and Federico Micheli. A Multicenter Comparative Study of Impulse Control Disorder in Latin American Patients With Parkinson Disease. *Clinical Neuropharmacology*, 40(2):51–55, April 2017. ISSN 1537-162X. doi: 10.1097/WNF.0000000000000202.
- O. Rascol, D. J. Brooks, A. D. Korczyn, P. P. De Deyn, C. E. Clarke, and A. E. Lang. A five-year study of the incidence of dyskinesia in patients with early Parkinson’s disease who were treated with ropinirole or levodopa. *The New England Journal of Medicine*, 342(20):1484–1491, May 2000. ISSN 0028-4793. doi: 10.1056/NEJM200005183422004.
- Sara Redenek, Duan Flisar, Maja Kojovi, Milica Gregori Kramberger, Dejan Georgiev, Zvezdan Pirtoek, Maja Trot, and Vita Dolan. Dopaminergic Pathway Genes Influence Adverse Events Related to Dopaminergic Treatment in Parkinson’s Disease. *Frontiers in Pharmacology*, 10:8, 2019. ISSN 1663-9812. doi: 10.3389/fphar.2019.00008.
- Matthew R. Robinson, Aaron Kleinman, Mariaelisa Graff, Anna A. E. Vinkhuyzen, David Couper, Michael B. Miller, Wouter J. Peyrot, Abdel Abdellaoui, Brendan P. Zietsch, Ilja M. Nolte, Jana V. van Vliet-Ostaptchouk, Harold Snieder, Sarah E. Medland, Nicholas G. Martin, Patrik K. E. Magnusson, William G. Iacono, Matt McGue, Kari E. North, Jian Yang, and Peter M. Visscher. Genetic evidence of assortative mating in humans. *Nature Human Behaviour*, 1(1):1–13, January 2017. ISSN 2397-3374. doi: 10.1038/s41562-016-0016. URL <https://www.nature.com/articles/s41562-016-0016>. Number: 1 Publisher: Nature Publishing Group.
- Mayela Rodríguez-Violante, Paulina González-Latapi, Amin Cervantes-Arriaga, Azyadeh Camacho-Ordoñez, and Daniel Weintraub. Impulse control and related disorders in Mexican Parkinson’s disease patients. *Parkinsonism & Related Disorders*, 20(8):907–910, August 2014. ISSN 1873-5126. doi: 10.1016/j.parkreldis.2014.05.014.
- David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Learning representations by back-propagating errors. *Nature*, 323(6088):533–536, October 1986.

- ISSN 1476-4687. doi: 10.1038/323533a0. URL <https://www.nature.com/articles/323533a0>. Number: 6088 Publisher: Nature Publishing Group.
- G. Rylander. Psychoses and the punding and choreiform syndromes in addiction to central stimulant drugs. *Psychiatry, Neurologia, Neurochirurgia*, 75(3):203–212, June 1972. ISSN 0033-2666.
- Jasjeet Sachdeva, Vijay Harbishettar, Michelle Barraclough, Kathryn McDonald, and Iracema Leroi. Clinical profile of compulsive sexual behaviour and paraphilia in Parkinson’s disease. *Journal of Parkinson’s Disease*, 4(4):665–670, 2014. ISSN 1877-718X. doi: 10.3233/JPD-140366.
- Rachel E. Salas, Richard P. Allen, Christopher J. Earley, and Charlene E. Gamaldo. Drug hoarding: a case of atypical dopamine dysregulation syndrome in a RLS patient. *Movement Disorders: Official Journal of the Movement Disorder Society*, 24(4):627–628, March 2009. ISSN 1531-8257. doi: 10.1002/mds.22443.
- Ali Samii, John G. Nutt, and Bruce R. Ransom. Parkinson’s disease. *The Lancet*, 363(9423):1783–1793, May 2004. ISSN 0140-6736, 1474-547X. doi: 10.1016/S0140-6736(04)16305-8. URL [https://www.thelancet.com/journals/lancet/article/PIIS0140-6736\(04\)16305-8/abstract](https://www.thelancet.com/journals/lancet/article/PIIS0140-6736(04)16305-8/abstract). Publisher: Elsevier.
- Sandra Sanchez-Roige, Pierre Fontanillas, Sarah L. Elson, Joshua C. Gray, Harriet de Wit, James MacKillop, and Abraham A. Palmer. Genome-Wide Association Studies of Impulsive Personality Traits (BIS-11 and UPPS-P) and Drug Experimentation in up to 22,861 Adult Research Participants Identify Loci in the CACNA1I and CADM2 genes. *The Journal of Neuroscience: The Official Journal of the Society for Neuroscience*, 39(13):2562–2572, 2019. ISSN 1529-2401. doi: 10.1523/JNEUROSCI.2662-18.2019.
- Pournamy Sarathchandran, Sheena Soman, Gangadhara Sarma, Syam Krishnan, and Asha Kishore. Impulse control disorders and related behaviors in Indian patients with Parkinson’s disease. *Movement Disorders*, 28(13):1901–1902, 2013. ISSN 1531-8257. doi: 10.1002/mds.25557. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/mds.25557>. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/mds.25557>.
- Jeanne E. Savage, Philip R. Jansen, Sven Stringer, Kyoko Watanabe, Julien Bryois, Christiaan A. de Leeuw, Mats Nagel, Swapnil Awasthi, Peter B. Barr, Jonathan R. I. Coleman, Katrina L. Grasby, Anke R. Hammerschlag, Jakob A. Kaminski, Robert Karlsson, Eva Krapohl, Max Lam, Marianne Nygaard, Chandra A. Reynolds, Joey W. Trampush, Hannah Young, Delilah Zabaneh, Sara Hägg, Narelle K. Hansell, Ida K. Karlsson, Sten Linnarsson, Grant W. Montgomery, Ana B. Muñoz-Manchado,

- Erin B. Quinlan, Gunter Schumann, Nathan G. Skene, Bradley T. Webb, Tonya White, Dan E. Arking, Dimitrios Avramopoulos, Robert M. Bilder, Panos Bitsios, Katherine E. Burdick, Tyrone D. Cannon, Ornit Chiba-Falek, Andrea Christoforou, Elizabeth T. Cirulli, Eliza Congdon, Aiden Corvin, Gail Davies, Ian J. Deary, Pamela DeRosse, Dwight Dickinson, Srdjan Djurovic, Gary Donohoe, Emily Drabant Conley, Johan G. Eriksson, Thomas Espeseth, Nelson A. Freimer, Stella Giakoumaki, Ina Giegling, Michael Gill, David C. Glahn, Ahmad R. Hariri, Alex Hatzimanolis, Matthew C. Keller, Emma Knowles, Deborah Koltai, Bettina Konte, Jari Lahti, Stephanie Le Hellard, Todd Lencz, David C. Liewald, Edythe London, Astri J. Lundervold, Anil K. Malhotra, Ingrid Melle, Derek Morris, Anna C. Need, William Ollier, Aarno Palotie, Antony Payton, Neil Pendleton, Russell A. Poldrack, Katri Räikkönen, Ivar Reinvang, Panos Roussos, Dan Rujescu, Fred W. Sabb, Matthew A. Scult, Olav B. Smeland, Nikolaos Smyrnis, John M. Starr, Vidar M. Steen, Nikos C. Stefanis, Richard E. Straub, Kjetil Sundet, Henning Tiemeier, Aristotle N. Voineskos, Daniel R. Weinberger, Elisabeth Widen, Jin Yu, Goncalo Abecasis, Ole A. Andreassen, Gerome Breen, Lene Christiansen, Birgit Debrabant, Danielle M. Dick, Andreas Heinz, Jens Hjerling-Leffler, M. Arfan Ikram, Kenneth S. Kendler, Nicholas G. Martin, Sarah E. Medland, Nancy L. Pedersen, Robert Plomin, Tinca J. C. Polderman, Stephan Ripke, Sophie van der Sluis, Patrick F. Sullivan, Scott I. Vrieze, Margaret J. Wright, and Danielle Posthuma. Genome-wide association meta-analysis in 269,867 individuals identifies new genetic and functional links to intelligence. *Nature Genetics*, 50(7): 912–919, 2018. ISSN 1546-1718. doi: 10.1038/s41588-018-0152-6.
- E. Schiørring. Psychopathology induced by "speed drugs". *Pharmacology, Biochemistry, and Behavior*, 14 Suppl 1:109–122, 1981. ISSN 0091-3057.
- S. Schlosser, D. W. Black, S. Repertinger, and D. Freet. Compulsive buying. Demography, phenomenology, and comorbidity in 46 subjects. *General Hospital Psychiatry*, 16(3):205–212, May 1994. ISSN 0163-8343. doi: 10.1016/0163-8343(94)90103-1.
- Michael K. Scullin, Ann B. Sollinger, Julia Land, Cathy Wood-Siverio, Lavezza Zanders, Raven Lee, Alan Freeman, Felicia C. Goldstein, Donald L. Bliwise, and Stewart A. Factor. Sleep and impulsivity in Parkinson's disease. *Parkinsonism & Related Disorders*, 19(11):991–994, November 2013. ISSN 1873-5126. doi: 10.1016/j.parkreldis.2013.06.018.
- Skipper Seabold and Josef Perktold. Statsmodels: Econometric and Statistical Modeling with Python. *Proceedings of the 9th Python in Science Conference*, 2010, January 2010.
- Tanya Simuni, Michael C. Brumm, Liz Uribe, Chelsea Caspell-Garcia, Christopher S. Coffey, Andrew Siderowf, Roy N. Alcalay, John Q. Trojanowski, Leslie M. Shaw, John

- Seibyl, Andrew Singleton, Arthur W. Toga, Doug Galasko, Tatiana Foroud, Kelly Nudelman, Duygu Tosun-Turgut, Kathleen Poston, Daniel Weintraub, Brit Mollenhauer, Caroline M. Tanner, Karl Kieburtz, Lana M. Chahine, Alyssa Reimer, Samantha Hutten, Susan Bressman, Kenneth Marek, and Parkinson's Progression Markers Initiative Investigators. Clinical and Dopamine Transporter Imaging Characteristics of Leucine Rich Repeat Kinase 2 (LRRK2) and Glucosylceramidase Beta (GBA) Parkinson's Disease Participants in the Parkinson's Progression Markers Initiative: A Cross-Sectional Study. *Movement Disorders: Official Journal of the Movement Disorder Society*, 35(5):833–844, 2020. ISSN 1531-8257. doi: 10.1002/mds.27989.
- Ashley H. Spencer, Hugh Rickards, Alfonso Fasano, and Andrea E. Cavanna. The prevalence and clinical characteristics of punding in Parkinson's disease. *Movement Disorders: Official Journal of the Movement Disorder Society*, 26(4):578–586, March 2011. ISSN 1531-8257. doi: 10.1002/mds.23508.
- Maria Grazia Spillantini, Marie Luise Schmidt, Virginia M.-Y. Lee, John Q. Trojanowski, Ross Jakes, and Michel Goedert. -Synuclein in Lewy bodies. *Nature*, 388(6645):839–840, August 1997. ISSN 1476-4687. doi: 10.1038/42166. URL <https://www.nature.com/articles/42166>. Number: 6645 Publisher: Nature Publishing Group.
- N. Sáez-Francàs, G. Martí Andrés, N. Ramírez, O. de Fàbregues, J. Álvarez Sabín, M. Casas, and J. Hernández-Vara. [Clinical and psychopathological factors associated with impulse control disorders in Parkinson's disease]. *Neurologia (Barcelona, Spain)*, 31(4):231–238, May 2016. ISSN 1578-1968. doi: 10.1016/j.nrl.2015.05.002.
- Lut Tamam, Mehtap Bican, and Necla Keskin. Impulse control disorders in elderly patients. *Comprehensive Psychiatry*, 55(4):1022–1028, May 2014. ISSN 1532-8384. doi: 10.1016/j.comppsy.2013.12.003.
- Robert Tibshirani. Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):267–288, 1996. ISSN 0035-9246. URL <https://www.jstor.org/stable/2346178>. Publisher: [Royal Statistical Society, Wiley].
- Andrey N. Tikhonov, Vasiliy Y. Arsenin, and Fritz John. *Solutions of Ill Posed Problems*. John Wiley & Sons Inc, Washington : New York, August 1977. ISBN 978-0-470-99124-4.
- Guilherme T. Valença, Philip G. Glass, Nadja N. Negreiros, Meirelayne B. Duarte, Lais M. G. B. Ventura, Mila Mueller, and Jamarly Oliveira-Filho. Past smoking and current dopamine agonist use show an independent and dose-dependent association with impulse control disorders in Parkinson's disease. *Parkinsonism & Related Disorders*, 19(7):698–700, July 2013. ISSN 1873-5126. doi: 10.1016/j.parkreldis.2013.03.004.

Annamaria Vallelunga, Raffaella Flaibani, Patrizia Formento-Dojot, Roberta Biundo, Silvia Facchini, and Angelo Antonini. Role of genetic polymorphisms of the dopaminergic system in Parkinson's disease patients with impulse control disorders. *Parkinsonism & Related Disorders*, 18(4):397–399, May 2012. ISSN 1873-5126. doi: 10.1016/j.parkreldis.2011.10.019.

Guy Van Camp, Anja Flamez, Bernard Cosyns, Caroline Weytjens, Luc Muyldermans, Michel Van Zandijcke, Johan De Sutter, Patrick Santens, Pierre Decoodt, Christian Moerman, and Danny Schoors. Treatment of Parkinson's disease with pergolide and relation to restrictive valvular heart disease. *Lancet (London, England)*, 363(9416): 1179–1183, April 2004. ISSN 1474-547X. doi: 10.1016/S0140-6736(04)15945-X.

Stéphanie M. van den Berg, Marleen H. M. de Moor, Karin J. H. Verweij, Robert F. Krueger, Michelle Luciano, Alejandro Arias Vasquez, Lindsay K. Matteson, Jaime Derringer, Tõnu Esko, Najaf Amin, Scott D. Gordon, Narelle K. Hansell, Amy B. Hart, Ilkka Seppälä, Jennifer E. Huffman, Bettina Konte, Jari Lahti, Minyoung Lee, Mike Miller, Teresa Natile, Toshiko Tanaka, Alexander Teumer, Alexander Viktorin, Juho Wedenoja, Abdel Abdellaoui, Goncalo R. Abecasis, Daniel E. Adkins, Arpana Agrawal, Jüri Allik, Katja Appel, Timothy B. Bigdeli, Fabio Busonero, Harry Campbell, Paul T. Costa, George Davey Smith, Gail Davies, Harriet de Wit, Jun Ding, Barbara E. Engelhardt, Johan G. Eriksson, Iryna O. Fedko, Luigi Ferrucci, Barbara Franke, Ina Giegling, Richard Grucza, Annette M. Hartmann, Andrew C. Heath, Kati Heinonen, Anjali K. Henders, Georg Homuth, Jouke-Jan Hottenga, William G. Iacono, Joost Janzing, Markus Jokela, Robert Karlsson, John P. Kemp, Matthew G. Kirkpatrick, Antti Latvala, Terho Lehtimäki, David C. Liewald, Pamela A. F. Madden, Chiara Magri, Patrik K. E. Magnusson, Jonathan Marten, Andrea Maschio, Hamdi Mbarek, Sarah E. Medland, Evelin Mihailov, Yuri Milaneschi, Grant W. Montgomery, Matthias Nauck, Michel G. Nivard, Klaasjan G. Ouwens, Aarno Palotie, Erik Pettersson, Ozren Polasek, Yong Qian, Laura Pulkki-Råback, Olli T. Raitakari, Anu Realo, Richard J. Rose, Daniela Ruggiero, Carsten O. Schmidt, Wendy S. Slutske, Rossella Sorice, John M. Starr, Beate St Pourcain, Angelina R. Sutin, Nicholas J. Timpson, Holly Trochet, Sita Vermeulen, Eero Vuoksima, Elisabeth Widen, Jasper Wouda, Margaret J. Wright, Lina Zgaga, Generation Scotland, David Porteous, Alessandra Minelli, Abraham A. Palmer, Dan Rujescu, Marina Ciullo, Caroline Hayward, Igor Rudan, Andres Metspalu, Jaakko Kaprio, Ian J. Deary, Katri Räikkönen, James F. Wilson, Liisa Keltikangas-Järvinen, Laura J. Bierut, John M. Hettema, Hans J. Grabe, Brenda W. J. H. Penninx, Cornelia M. van Duijn, David M. Evans, David Schlessinger, Nancy L. Pedersen, Antonio Terracciano, Matt McGue, Nicholas G. Martin, and Dorret I. Boomsma. Meta-analysis of Genome-Wide Association Studies for Extraversion: Findings from the Genetics of Personality Con-

- sortium. *Behavior Genetics*, 46(2):170–182, March 2016. ISSN 1573-3297. doi: 10.1007/s10519-015-9735-5.
- Guido Van Rossum and Fred L. Drake. *Python 3 Reference Manual*. CreateSpace, Scotts Valley, CA, 2009. ISBN 978-1-4414-1269-0.
- Vladimir N. Vapnik and Alexander Lerner. Pattern recognition using generalized portrait method. *Automation and Remote Control*, 24:774–780, 1963.
- L. Vela, J. C. Martínez Castrillo, P. García Ruiz, C. Gasca-Salas, Y. Macías Macías, E. Pérez Fernández, I. Ybot, E. Lopez Valdés, M. M. Kurtis, I. J. Posada Rodriguez, M. Mata, C. Ruiz Huete, M. Eimil, C. Borrue, J. Del Val, L. López-Manzanares, A. Rojo Sebastian, and R. Marasescu. The high prevalence of impulse control behaviors in patients with early-onset Parkinson’s disease: A cross-sectional multicenter study. *Journal of the Neurological Sciences*, 368:150–154, September 2016. ISSN 1878-5883. doi: 10.1016/j.jns.2016.07.003.
- Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C. J. Carey, Ihan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, and Paul van Mulbregt. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nature Methods*, 17(3):261–272, March 2020. ISSN 1548-7105. doi: 10.1038/s41592-019-0686-2. URL <https://www.nature.com/articles/s41592-019-0686-2>. Number: 3 Publisher: Nature Publishing Group.
- Martine Visser, Dagmar Verbaan, Stephanie M. van Rooden, Anne M. Stiggelbout, Johan Marinus, and Jacobus J. van Hilten. Assessment of psychiatric complications in Parkinson’s disease: The SCOPA-PC. *Movement Disorders: Official Journal of the Movement Disorder Society*, 22(15):2221–2228, November 2007. ISSN 0885-3185. doi: 10.1002/mds.21696.
- Valerie Voon, Mandy Sohr, Anthony E. Lang, Marc N. Potenza, Andrew D. Siderowf, Jacqueline Whetteckey, Daniel Weintraub, Glen R. Wunderlich, and Mark Stacy. Impulse control disorders in Parkinson disease: a multicenter case-control study. *Annals of Neurology*, 69(6):986–996, June 2011. ISSN 1531-8249. doi: 10.1002/ana.22356.
- Ying Wang, Jing Guo, Guiyan Ni, Jian Yang, Peter M. Visscher, and Loic Yengo. Theoretical and empirical quantification of the accuracy of polygenic scores in ancestry

divergent populations. *Nature Communications*, 11(1):3865, July 2020. ISSN 2041-1723. doi: 10.1038/s41467-020-17719-y. URL <https://www.nature.com/articles/s41467-020-17719-y>. Number: 1 Publisher: Nature Publishing Group.

Mathilde Wanneveich, Frédéric Moisan, Hélène Jacqmin-Gadda, Alexis Elbaz, and Pierre Joly. Projections of prevalence, lifetime risk, and life expectancy of Parkinson’s disease (2010-2030) in France. *Movement Disorders: Official Journal of the Movement Disorder Society*, 33(9):1449–1455, 2018. ISSN 1531-8257. doi: 10.1002/mds.27447.

C. D. Ward and W. R. Gibb. Research diagnostic criteria for Parkinson’s disease. *Advances in Neurology*, 53:245–249, 1990. ISSN 0091-3952.

Hunna J. Watson, Zeynep Yilmaz, Laura M. Thornton, Christopher Hübel, Jonathan R. I. Coleman, Hélène A. Gaspar, Julien Bryois, Anke Hinney, Virpi M. Leppä, Manuel Mattheisen, Sarah E. Medland, Stephan Ripke, Shuyang Yao, Paola Giusti-Rodríguez, Anorexia Nervosa Genetics Initiative, Ken B. Hanscombe, Kirstin L. Purves, Eating Disorders Working Group of the Psychiatric Genomics Consortium, Roger A. H. Adan, Lars Alfredsson, Tetsuya Ando, Ole A. Andreassen, Jessica H. Baker, Wade H. Berrettini, Ilka Boehm, Claudette Boni, Vesna Boraska Perica, Katharina Buehren, Roland Burghardt, Matteo Cassina, Sven Cichon, Maurizio Clementi, Roger D. Cone, Philippe Courtet, Scott Crow, James J. Crowley, Unna N. Danner, Oliver S. P. Davis, Martina de Zwaan, George Dedoussis, Daniela Degortes, Janiece E. DeSocio, Danielle M. Dick, Dimitris Dikeos, Christian Dina, Monika Dmitrzak-Weglarz, Elisa Docampo, Laramie E. Duncan, Karin Egberts, Stefan Ehrlich, Geòrgia Escaramís, Tõnu Esko, Xavier Estivill, Anne Farmer, Angela Favaro, Fernando Fernández-Aranda, Manfred M. Fichter, Krista Fischer, Manuel Föcker, Lenka Foretova, Andreas J. Forstner, Monica Forzan, Christopher S. Franklin, Steven Gallinger, Ina Giegling, Johanna Giuranna, Fragiskos Gonidakis, Philip Gormwood, Monica Gratacos Mayora, Sébastien Guillaume, Yiran Guo, Hakon Hakonarson, Konstantinos Hatzikotoulas, Joanna Hauser, Johannes Hebebrand, Sietske G. Helder, Stefan Herms, Beate Herpertz-Dahlmann, Wolfgang Herzog, Laura M. Huckins, James I. Hudson, Hartmut Imgart, Hidetoshi Inoko, Vladimir Janout, Susana Jiménez-Murcia, Antonio Julià, Gursharan Kalsi, Deborah Kaminská, Jaakko Kaprio, Leila Karhunen, Andreas Karwautz, Martien J. H. Kas, James L. Kennedy, Anna Keski-Rahkonen, Kirsty Kiezebrink, Youl-Ri Kim, Lars Klareskog, Kelly L. Klump, Gun Peggy S. Knudsen, Maria C. La Via, Stephanie Le Hellard, Robert D. Levitan, Dong Li, Lisa Lilienfeld, Bochao Danae Lin, Jolanta Lissowska, Jurjen Luykx, Pierre J. Magistretti, Mario Maj, Katrin Mannik, Sara Marsal, Christian R. Marshall, Morten Mattingsdal, Sara McDevitt, Peter McGuffin, Andres Metspalu, Ingrid Meulenbelt, Nadia Micali, Karen Mitchell, Alessio Maria Monteleone, Palmiero Monteleone, Melissa A. Munn-Chernoff, Benedetta Nacmias, Marie Navratilova,

- Ioanna Ntalla, Julie K. O'Toole, Roel A. Ophoff, Leonid Padyukov, Aarno Palotie, Jacques Pantel, Hana Papezova, Dalila Pinto, Raquel Rabionet, Anu Raevuori, Nicolas Ramoz, Ted Reichborn-Kjennerud, Valdo Ricca, Samuli Ripatti, Franziska Ritschel, Marion Roberts, Alessandro Rotondo, Dan Rujescu, Filip Rybakowski, Paolo Santonastaso, André Scherag, Stephen W. Scherer, Ulrike Schmidt, Nicholas J. Schork, Alexandra Schosser, Jochen Seitz, Lenka Slachtova, P. Eline Slagboom, Margarita C. T. Slof-Op 't Landt, Agnieszka Slopian, Sandro Sorbi, Beata witkowska, Jin P. Szatkiewicz, Ioanna Tachmazidou, Elena Tenconi, Alfonso Tortorella, Federica Tozzi, Janet Treasure, Artemis Tsitsika, Marta Tyszkiewicz-Nwafor, Konstantinos Tziouvas, Annemarie A. van Elburg, Eric F. van Furth, Gudrun Wagner, Esther Walton, Elisabeth Widen, Eleftheria Zeggini, Stephanie Zerwas, Stephan Zipfel, Andrew W. Bergen, Joseph M. Boden, Harry Brandt, Steven Crawford, Katherine A. Halmi, L. John Horwood, Craig Johnson, Allan S. Kaplan, Walter H. Kaye, James E. Mitchell, Catherine M. Olsen, John F. Pearson, Nancy L. Pedersen, Michael Strober, Thomas Werge, David C. Whiteman, D. Blake Woodside, Garret D. Stuber, Scott Gordon, Jakob Grove, Anjali K. Henders, Anders Juréus, Katherine M. Kirk, Janne T. Larsen, Richard Parker, Liselotte Petersen, Jennifer Jordan, Martin Kennedy, Grant W. Montgomery, Tracey D. Wade, Andreas Birgegård, Paul Lichtenstein, Claes Norring, Mikael Landén, Nicholas G. Martin, Preben Bo Mortensen, Patrick F. Sullivan, Gerome Breen, and Cynthia M. Bulik. Genome-wide association study identifies eight risk loci and implicates metabo-psychiatric origins for anorexia nervosa. *Nature Genetics*, 51(8):1207–1214, 2019. ISSN 1546-1718. doi: 10.1038/s41588-019-0439-2.
- Daniel Weintraub and Daniel O. Claassen. Impulse Control and Related Disorders in Parkinson's Disease. *International Review of Neurobiology*, 133:679–717, 2017. ISSN 2162-5514. doi: 10.1016/bs.irn.2017.04.006.
- Daniel Weintraub, Staci Hoops, Judy A. Shea, Kelly E. Lyons, Rajesh Pahwa, Erika D. Driver-Dunckley, Charles H. Adler, Marc N. Potenza, Janis Miyasaki, Andrew D. Siderowf, John E. Duda, Howard I. Hurtig, Amy Colcher, Stacy S. Horn, Matthew B. Stern, and Valerie Voon. Validation of the questionnaire for impulsive-compulsive disorders in Parkinson's disease. *Movement Disorders: Official Journal of the Movement Disorder Society*, 24(10):1461–1467, July 2009. ISSN 1531-8257. doi: 10.1002/mds.22571.
- Daniel Weintraub, Juergen Koester, Marc N. Potenza, Andrew D. Siderowf, Mark Stacy, Valerie Voon, Jacqueline Whetteckey, Glen R. Wunderlich, and Anthony E. Lang. Impulse control disorders in Parkinson disease: a cross-sectional study of 3090 patients. *Archives of Neurology*, 67(5):589–595, May 2010a. ISSN 1538-3687. doi: 10.1001/archneurol.2010.65.

- Daniel Weintraub, Mandy Sohr, Marc N. Potenza, Andrew D. Siderowf, Mark Stacy, Valerie Voon, Jacqueline Whetteckey, Glen R. Wunderlich, and Anthony E. Lang. Amantadine use associated with impulse control disorders in Parkinson disease in cross-sectional study. *Annals of Neurology*, 68(6):963–968, December 2010b. ISSN 1531-8249. doi: 10.1002/ana.22164.
- Daniel Weintraub, Eugenia Mamikonyan, Kimberly Papay, Judith A. Shea, Sharon X. Xie, and Andrew Siderowf. Questionnaire for Impulsive-Compulsive Disorders in Parkinson’s Disease-Rating Scale. *Movement Disorders: Official Journal of the Movement Disorder Society*, 27(2):242–247, February 2012. ISSN 1531-8257. doi: 10.1002/mds.24023.
- Naomi R. Wray, Michael E. Goddard, and Peter M. Visscher. Prediction of individual genetic risk to disease from genome-wide association studies. *Genome Research*, 17(10):1520–1528, October 2007. ISSN 1088-9051. doi: 10.1101/gr.6665407.
- Jian Yang, S. Hong Lee, Michael E. Goddard, and Peter M. Visscher. GCTA: a tool for genome-wide complex trait analysis. *American Journal of Human Genetics*, 88(1):76–82, January 2011. ISSN 1537-6605. doi: 10.1016/j.ajhg.2010.11.011.
- Zheng Ye, Anke Hammer, Estela Camara, and Thomas F. Münte. Pramipexole modulates the neural network of reward anticipation. *Human Brain Mapping*, 32(5):800–811, May 2011. ISSN 1097-0193. doi: 10.1002/hbm.21067.
- Loic Yengo, Julia Sidorenko, Kathryn E. Kemper, Zhili Zheng, Andrew R. Wood, Michael N. Weedon, Timothy M. Frayling, Joel Hirschhorn, Jian Yang, Peter M. Visscher, and GIANT Consortium. Meta-analysis of genome-wide association studies for height and body mass index in 700000 individuals of European ancestry. *Human Molecular Genetics*, 27(20):3641–3649, 2018. ISSN 1460-2083. doi: 10.1093/hmg/ddy271.
- Hana You, Louise-Laure Mariani, Graziella Mangone, Delphine Le Febvre de Nailly, Fanny Charbonnier-Beaupel, and Jean-Christophe Corvol. Molecular basis of dopamine replacement therapy and its side effects in Parkinson’s disease. *Cell and Tissue Research*, 373(1):111–135, 2018. ISSN 1432-0878. doi: 10.1007/s00441-018-2813-2.
- Kimberly S. Young. Internet Addiction: The Emergence of a New Clinical Disorder. *CyberPsychology & Behavior*, 1(3):237–244, January 1998. ISSN 1094-9313. doi: 10.1089/cpb.1998.1.237. URL <https://www.liebertpub.com/doi/10.1089/cpb.1998.1.237>. Publisher: Mary Ann Liebert, Inc., publishers.
- Insha Zahoor, Amrina Shafi, and Ehtishamul Haq. Pharmacological Treatment of Parkinsons Disease. In Thomas B. Stoker and Julia C. Greenland, editors, *Parkin-*

sons Disease: Pathogenesis and Clinical Aspects. Codon Publications, Brisbane (AU), 2018. ISBN 978-0-9944381-6-4. URL <http://www.ncbi.nlm.nih.gov/books/NBK536726/>.

Shahidee Zainal Abidin, Eng Liang Tan, Soon-Choy Chan, Ameerah Jaafar, Alex Xuen Lee, Mohd Hamdi Noor Abd Hamid, Nor Azian Abdul Murad, Nur Fadlina Pakarul Razy, Shahrul Azmin, Azlina Ahmad Annuar, Shen Yang Lim, Pike-See Cheah, King-Hwa Ling, and Norlinah Mohamed Ibrahim. DRD and GRIN2B polymorphisms and their association with the development of impulse control behaviour among Malaysian Parkinson's disease patients. *BMC neurology*, 15:59, April 2015. ISSN 1471-2377. doi: 10.1186/s12883-015-0316-2.

Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005. ISSN 1467-9868. doi: 10.1111/j.1467-9868.2005.00503.x. URL <https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/j.1467-9868.2005.00503.x>. _eprint: <https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/j.1467-9868.2005.00503.x>.